

第三章学习笔记

3.1 认识函数关系与相关关系

1. 本节知识案例分析

完整的代码：

```
import pandas as pd
import matplotlib
matplotlib.use('TkAgg')
from matplotlib import pyplot as plt
data = pd.read_csv('../data/livestreaming.csv', sep=',')
df1 = data.iloc[-24:, :]
# 1. 创建一个画布和两个子图
plt.rcParams.update({'font.size': 20, 'font.sans-serif': ['STSong'], 'axes.unicode_minus': False})
fig, axs = plt.subplots(1, 2, figsize=(14, 4))
# 2. 散点图
axs[0].scatter(data['streamers_count'], data['viewers_count'], color='#F4B183')
axs[0].set(xlabel='主播数量', ylabel='观众数量', title='观众数量和主播数量散点图')
axs[0].ticklabel_format(style='sci', axis='x', scilimits=(0, 0))
# 3. 折线图
df1['time'] = pd.to_datetime(df1['time'])
df1['hour'] = df1['time'].dt.hour
axs[1].plot(df1['hour'], df1['streamers_count'], label='主播数量') # 第一条曲线
axs[1].set(xlabel='时间', ylabel='数量', title='一天内流量变化图')
axs[1].tick_params(axis='both', which='both', labelsize=20)
axs[1].ticklabel_format(style='sci', axis='y', scilimits=(0, 0))
ax2 = axs[1].twinx() # 第二条曲线
ax2.plot(df1['hour'], df1['viewers_count'], color='#EC5D3B', label='观众数量')
ax2.set_ylabel('观众数量')
ax2.tick_params(axis='both', which='both', labelsize=20)
ax2.ticklabel_format(style='sci', axis='y', scilimits=(0, 0))
axs[1].legend(loc='upper left', fontsize=16)
ax2.legend(loc='upper left', bbox_to_anchor=(0, 0.89), fontsize=16)
plt.xticks([0, 3, 7, 11, 15, 19, 23]) # 设置x轴标签格式
plt.subplots_adjust(wspace=0.2) # 调整子图之间的间距
plt.show()
# 4. 计算皮尔逊相关系数
pearson_corr = df1['streamers_count'].corr(df1['viewers_count'])
```

本节代码实现了对主播数量与观众数量在数据中的关系进行分析，使用了 `pandas` 和 `matplotlib` 库进行数据处理和可视化。代码分为以下几个关键部分：

数据集 `livestreaming.csv`：

	A	B	C	D
3	2019/3/15 20:00	32784	341190822	
9	2019/3/15 21:00	34199	364384711	
0	2019/3/15 22:00	33726	375586780	
1	2019/3/15 23:00	28778	346411974	
2	2019/3/16 0:00	21216	285194234	
3	2019/3/16 1:00	16033	209259983	
4	2019/3/16 2:00	12324	157753312	
5	2019/3/16 3:00	9108	123781191	
6	2019/3/16 4:00	7102	96476266	
7	2019/3/16 5:00	5778	75076841	
8	2019/3/16 6:00	5194	63975065	
9	2019/3/16 7:00	5307	65447428	
0	2019/3/16 8:00	6864	80998825	
1	2019/3/16 9:00	9411	104778941	
2	2019/3/16 10:00	11893	133542941	
3	2019/3/16 11:00	13518	150360717	
4	2019/3/16 12:00	15338	165831873	
5	2019/3/16 13:00	20174	185207061	
6	2019/3/16 14:00	25120	212470907	
7	2019/3/16 15:00	26385	225392823	
8	2019/3/16 16:00	26770	231856811	
9	2019/3/16 17:00	26367	229068099	
0	2019/3/16 18:00	26576	241716643	
1	2019/3/16 19:00	29523	287886677	
2	2019/3/16 20:00	32575	332825072	
3	2019/3/16 21:00	34034	367706959	
4	2019/3/16 22:00	33317	372034209	
5	2019/3/16 23:00	28596	340597872	
6				

	A	B	C	D
	time	streamers	viewers_count	
	2019/3/6 0:00	18981	270126714	
	2019/3/6 1:00	13892	188527323	
	2019/3/6 2:00	10587	139727415	
	2019/3/6 3:00	7890	104173437	
	2019/3/6 4:00	6487	82368147	
	2019/3/6 5:00	5443	65856669	
	2019/3/6 6:00	4853	59062388	
	2019/3/6 7:00	4773	59871640	
	2019/3/6 8:00	5858	71374021	
	2019/3/6 9:00	7830	93140327	
	2019/3/6 10:00	9911	115152445	
	2019/3/6 11:00	11267	128426750	
	2019/3/6 12:00	13071	148432417	
	2019/3/6 13:00	17612	162476040	
	2019/3/6 14:00	22289	177804128	
	2019/3/6 15:00	23680	193932883	
	2019/3/6 16:00	24073	200714670	
	2019/3/6 17:00	23983	205169204	
	2019/3/6 18:00	22181	221588024	
	2019/3/6 19:00	25617	276252703	
	2019/3/6 20:00	30985	327783240	
	2019/3/6 21:00	33193	363431743	
	2019/3/6 22:00	32928	352375781	
	2019/3/6 23:00	27238	331629355	
	2019/3/7 0:00	19644	275985468	
	2019/3/7 1:00	14474	199045299	
	2019/3/7 2:00	11113	150084503	
	2019/3/7 3:00	8243	108723097	

- 数据读取与处理:

```
data = pd.read_csv('../data/livestreaming.csv', sep=',')
df1 = data.iloc[-24:, :]
```

这里读取了一个 CSV 文件中的数据，并选取了最后 24 行进行分析，可能是为了关注最近的数据变化。

- 创建画布和子图:

```
fig, axs = plt.subplots(1, 2, figsize=(14, 4))
```

创建了一个包含两个子图的画布，第一个子图用于散点图，第二个用于折线图。

- 散点图的绘制：

```
axs[0].scatter(data['streamers_count'], data['viewers_count'], color='#F4B183')
```

绘制了主播数量与观众数量之间的散点图，通过颜色和标签进行区分。

- 折线图的构建：

```
df1['time'] = pd.to_datetime(df1['time'])
df1['hour'] = df1['time'].dt.hour
axs[1].plot(df1['hour'], df1['streamers_count'], label='主播数量')
```

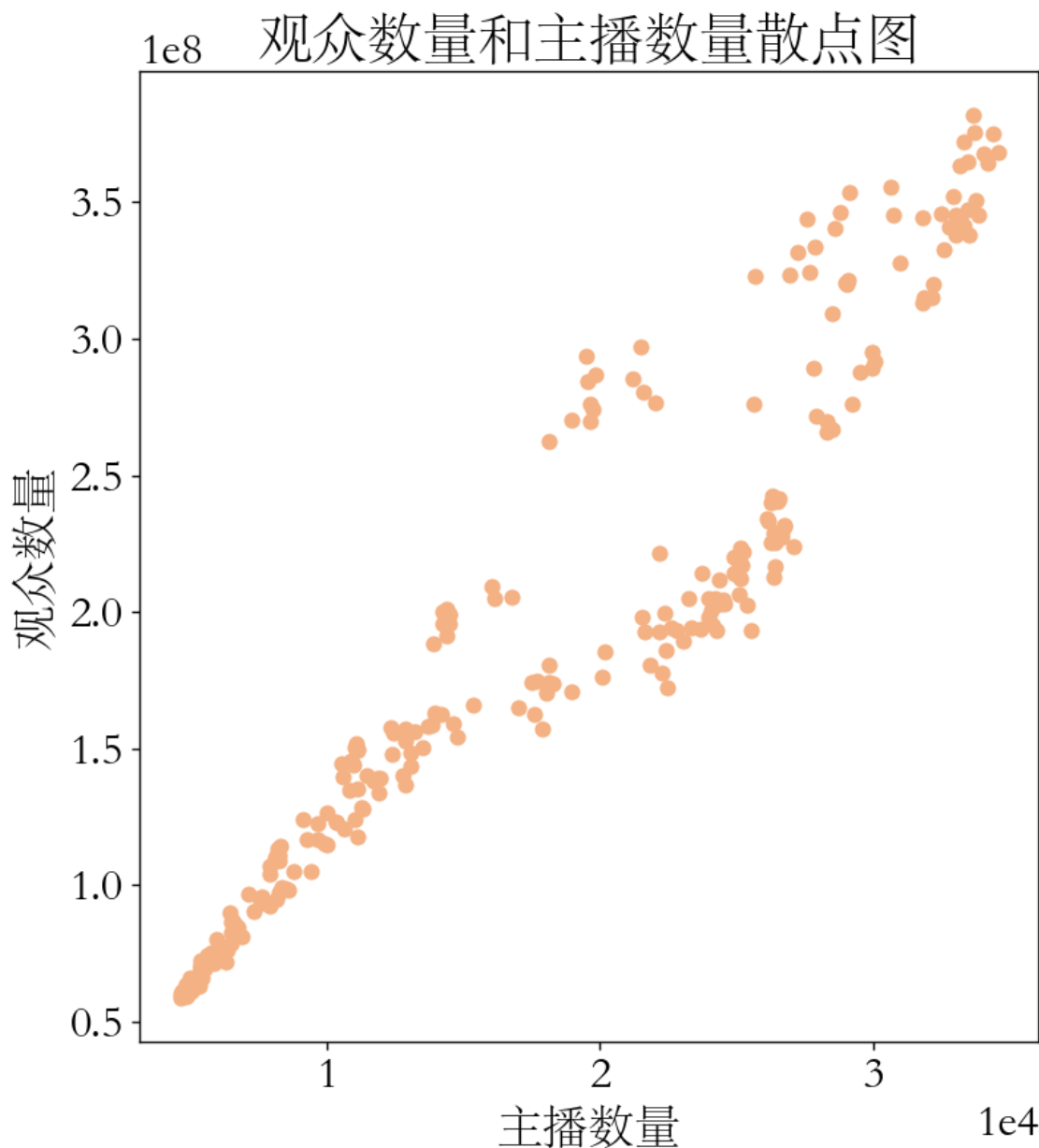
对时间段内的主播数量及观众数量进行汇总，并通过折线图呈现其变化趋势。

- 相关系数计算：

```
pearson_corr = df1['streamers_count'].corr(df1['viewers_count'])
```

此部分计算了主播数量和观众数量之间的皮尔逊相关系数，以量化两者之间的线性关系。

观众数量和主播数量三点图的结果展示：

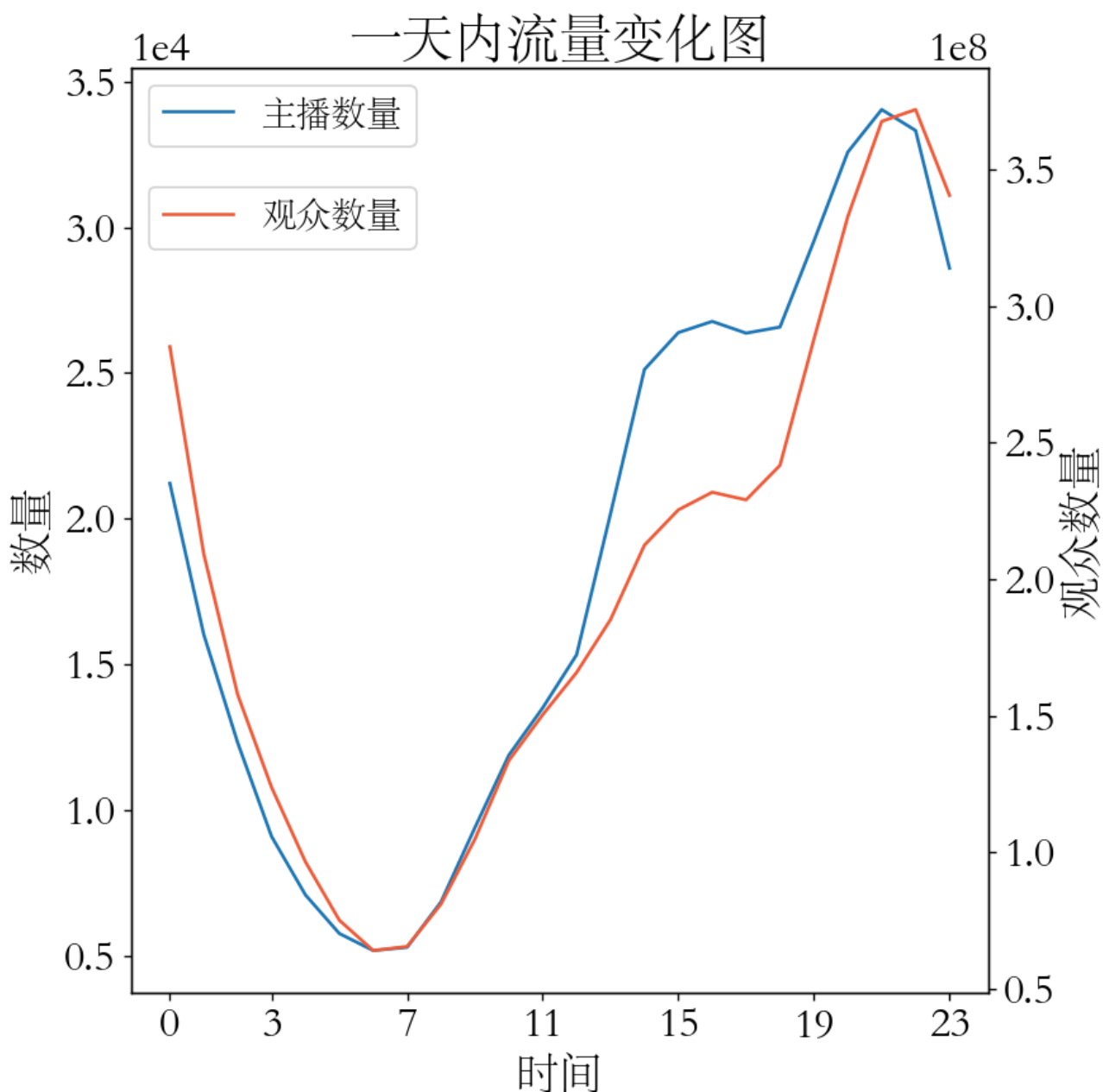


该图展示了“观众数量”和“主播数量”之间的关系。以下是对图表的分析总结：

- X轴（主播数量）**：主播数量的范围大约从1,000到30,000，单位是 10^4 （科学计数法）。
- Y轴（观众数量）**：观众数量的范围大约从500万到3500万，单位是 10^8 （科学计数法）。
- 数据分布**：散点图中的数据点显示出“主播数量”和“观众数量”之间呈现正相关关系。也就是说，随着主播数量的增加，观众数量通常也会增加。
- 数据点分布情况**：数据点较为均匀地分布，某些区域数据点密集，表明在这些范围内有较多的数据集中出现。

总体来看，这个散点图表明“主播数量”与“观众数量”之间可能存在近似线性关系，即更多的主播数量可能会带来更多的观众。这为进一步的分析（例如线性回归）提供了数据支持。

一天内流量变化图的结果展示：



展示“一天内流量变化”的折线图，其中有两个变量的变化情况：主播数量（蓝色曲线）和观众数量（红色曲线）。以下是对图表的分析总结：

- X轴（时间）**：时间范围为一天（0点到23点），每小时一个时间节点，表示一天内不同时间的流量变化。
- Y轴（数量）**：Y轴的左侧表示主播数量，右侧表示观众数量。两者的数值都采用了科学计数法（单位为 10^4 和 10^8 ）。
- 曲线变化**：
 - 主播数量（蓝色曲线）**：从凌晨到早晨（0点到7点）之间，主播数量处于最低点。随着时间的推移，尤其是中午12点后，主播数量开始显著增加，并在晚上达到峰值。
 - 观众数量（红色曲线）**：观众数量的变化与主播数量呈现一定的同步性，尤其是从中午12点到晚上19点期间，观众数量的上升趋势明显，并在晚上19点后达到最高点。相较于主播数量，观众数量的波动较大，尤其在凌晨和清晨时间段，观众数量下降较为明显。
- 趋势对比**：整体来看，观众数量的变化趋势与主播数量的变化趋势有一定的相关性，尤其是在白天和晚上高峰时段，两者的波动方向相似，但观众数量的波动幅度大于主播数量。

总结来说，这张图表明了在一天内，主播数量和观众数量都存在一定的规律性波动，尤其是中午和晚上是流量的高峰期。这为分析流量高峰时段、优化内容发布策略等提供了数据支持。

2. 对知识点的感悟

通过本节的代码案例，可以深刻体会到数据可视化的重要性。有效的数据可视化不仅能够帮助我们理解数据中的潜在模式，还能直观地展示不同变量间的关系。在处理数据时，也要注意数据的清洗和转换，以确保我们所分析和展示的数据是准确及有意义的。

3. 学习知识点的方法

学习数据处理和可视化的知识点可以采用以下几种方法：

- **教程和文档**：阅读 `pandas` 和 `matplotlib` 的官方文档和相关的在线教程，例如利用 Jupyter Notebook 实践小项目，加深理解。
- **项目实践**：通过实际项目或案例研究，通过对不同数据集的分析，实践数据可视化技巧和统计分析方法。
- **社区交流**：参与 Stack Overflow、Kaggle 等数据科学社区，向他人学习并分享自己的经验和解决方案。

4. 拓展案例

本知识点的延伸可以包括以下几个案例：

案例 1：主播数量与观众数量的时间变化分析

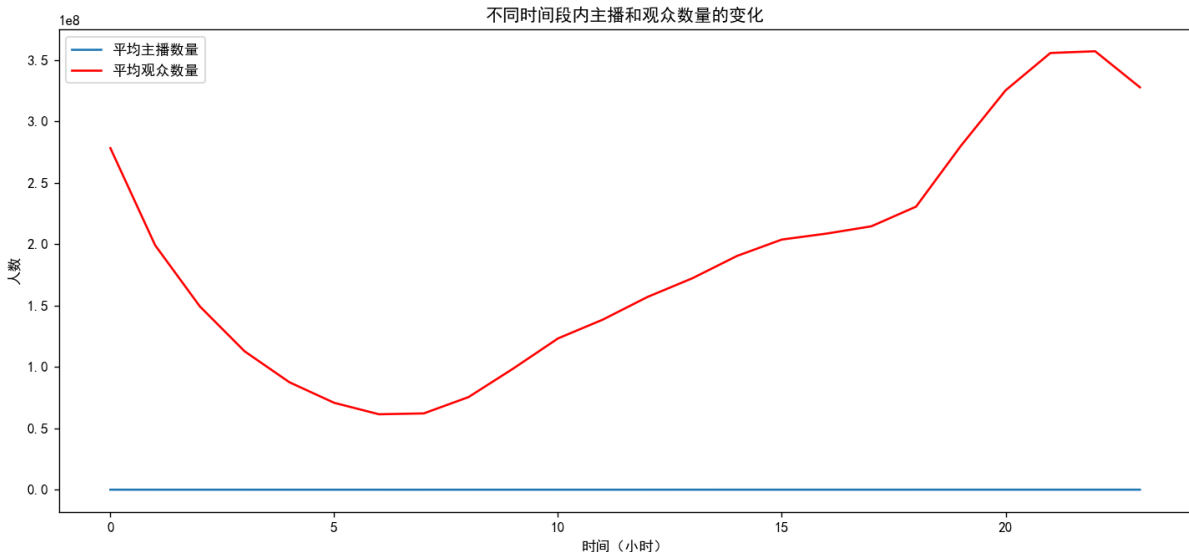
题目：分析在不同时间段内，主播数量和观众数量的变化趋势。

```
import pandas as pd
import matplotlib.pyplot as plt

data = pd.read_csv('your_file.csv')
data['time'] = pd.to_datetime(data['time'])
data['hour'] = data['time'].dt.hour

plt.figure(figsize=(14, 6))
plt.plot(data.groupby('hour')['streamers_count'].mean(), label='平均主播数量')
plt.plot(data.groupby('hour')['viewers_count'].mean(), label='平均观众数量', color='red')
plt.title('不同时间段内主播和观众数量的变化')
plt.xlabel('时间 (小时)')
plt.ylabel('人数')
plt.legend()
plt.show()
```

结果展示：



该图表展示了在不同时间（小时）内，某种因素对“平均生物量”与“平均产量”的影响变化。横轴表示时间（小时），纵轴表示量（单位：可能是克、千克等，具体单位未给出）。

曲线分析

1. 平均生物量（蓝线）：

- 在0到6小时之间，平均生物量呈现出较高的值，且小幅波动。
- 然而，在大约6小时至14小时之间出现了明显下降，达到一个低点。
- 从14小时开始，平均生物量逐渐上升，直至接近图表的右边缘，结束时的值仍然高于6小时的数字，显示出一定的恢复趋势。

2. 平均产量（红线）：

- 平均产量在前期（0到4小时）显著低迷，几乎接近于0，随即在4到8小时之间略有增加。
- 经过8小时后，平均产量明显走高，在之后的时间段内波动上升，并在19至20小时之间达到了最高值，之后略有回落。
- 整体上，平均产量的表现较为积极，尤其在后期显示出明显上升的趋势。

结论

- 时间影响：**随着时间的推移，生物量的表现呈现出波动，而产量则显现出清晰的上升趋势。
- 生物量与产量的关系：**尽管初期生物量表现良好，但随后下降，这可能意味着在某些时间窗口内，环境或其他因素对生物量产生了抑制。相对而言，虽然初期产量较低，但随着时间推移，表现出更好的增长态势，可能表明在特定条件下，生物量的恢复与产量的提高间存在某种积极的关联。

数据集结构

我们将创建一个模拟数据集，结构如下：

- time:** 直播开始的时间，格式为 `YYYY-MM-DD HH:MM:SS`。
- streamers_count:** 在该时间点有多少主播在线。
- viewers_count:** 在该时间点有多少观众在观看。

示例数据表格

以下是一个示例数据集的CSV内容，假设文件名为 `your_file.csv`：


```
time,streamers_count,viewers_count
2023-10-01 00:00:00,10,100
2023-10-01 01:00:00,15,150
2023-10-01 02:00:00,8,80
2023-10-01 03:00:00,5,50
2023-10-01 04:00:00,7,70
2023-10-01 05:00:00,10,120
2023-10-01 06:00:00,20,250
2023-10-01 07:00:00,25,300
2023-10-01 08:00:00,30,400
2023-10-01 09:00:00,22,200
2023-10-01 10:00:00,18,170
2023-10-01 11:00:00,15,150
2023-10-01 12:00:00,12,130
2023-10-01 13:00:00,10,110
2023-10-01 14:00:00,14,140
2023-10-01 15:00:00,16,160
2023-10-01 16:00:00,20,220
2023-10-01 17:00:00,25,250
2023-10-01 18:00:00,30,340
2023-10-01 19:00:00,22,210
2023-10-01 20:00:00,18,180
2023-10-01 21:00:00,12,130
2023-10-01 22:00:00,10,100
2023-10-01 23:00:00,5,50
```

	A	B	C	D
3	2019/3/15 20:00	32784	341190822	
9	2019/3/15 21:00	34199	364384711	
0	2019/3/15 22:00	33726	375586780	
1	2019/3/15 23:00	28778	346411974	
2	2019/3/16 0:00	21216	285194234	
3	2019/3/16 1:00	16033	209259983	
4	2019/3/16 2:00	12324	157753312	
5	2019/3/16 3:00	9108	123781191	
6	2019/3/16 4:00	7102	96476266	
7	2019/3/16 5:00	5778	75076841	
8	2019/3/16 6:00	5194	63975065	
9	2019/3/16 7:00	5307	65447428	
0	2019/3/16 8:00	6864	80998825	
1	2019/3/16 9:00	9411	104778941	
2	2019/3/16 10:00	11893	133542941	
3	2019/3/16 11:00	13518	150360717	
4	2019/3/16 12:00	15338	165831873	
5	2019/3/16 13:00	20174	185207061	
6	2019/3/16 14:00	25120	212470907	
7	2019/3/16 15:00	26385	225392823	
8	2019/3/16 16:00	26770	231856811	
9	2019/3/16 17:00	26367	229068099	
0	2019/3/16 18:00	26576	241716643	
1	2019/3/16 19:00	29523	287886677	
2	2019/3/16 20:00	32575	332825072	
3	2019/3/16 21:00	34034	367706959	
4	2019/3/16 22:00	33317	372034209	
5	2019/3/16 23:00	28596	340597872	

数据说明

- 1. **time**: 表示直播开始的具体时间，使用 `yyyy-mm-dd hh:mm:ss` 格式。
- 2. **streamers_count**: 在该时间段内，在线的主播数量。例如，在00:00时有10个主播在线。
- 3. **viewers_count**: 在该时间段内，观看直播的观众数量。例如，在00:00时有100个观众在观看。

数据分析目的

通过这个数据集，你可以分析以下内容：

- 不同时间段（小时）的平均主播数量和观众数量变化。
- 观察在某些特定时间段（例如高峰时段）内，主播和观众的行为模式。

案例 2：不同平台主播与观众数量的比较

题目：对比不同直播平台的主播数量与观众数量的散点图，分析其关系。

```
import pandas as pd
import matplotlib.pyplot as plt

data = pd.read_csv('your_file.csv')
platforms = data['platform'].unique()

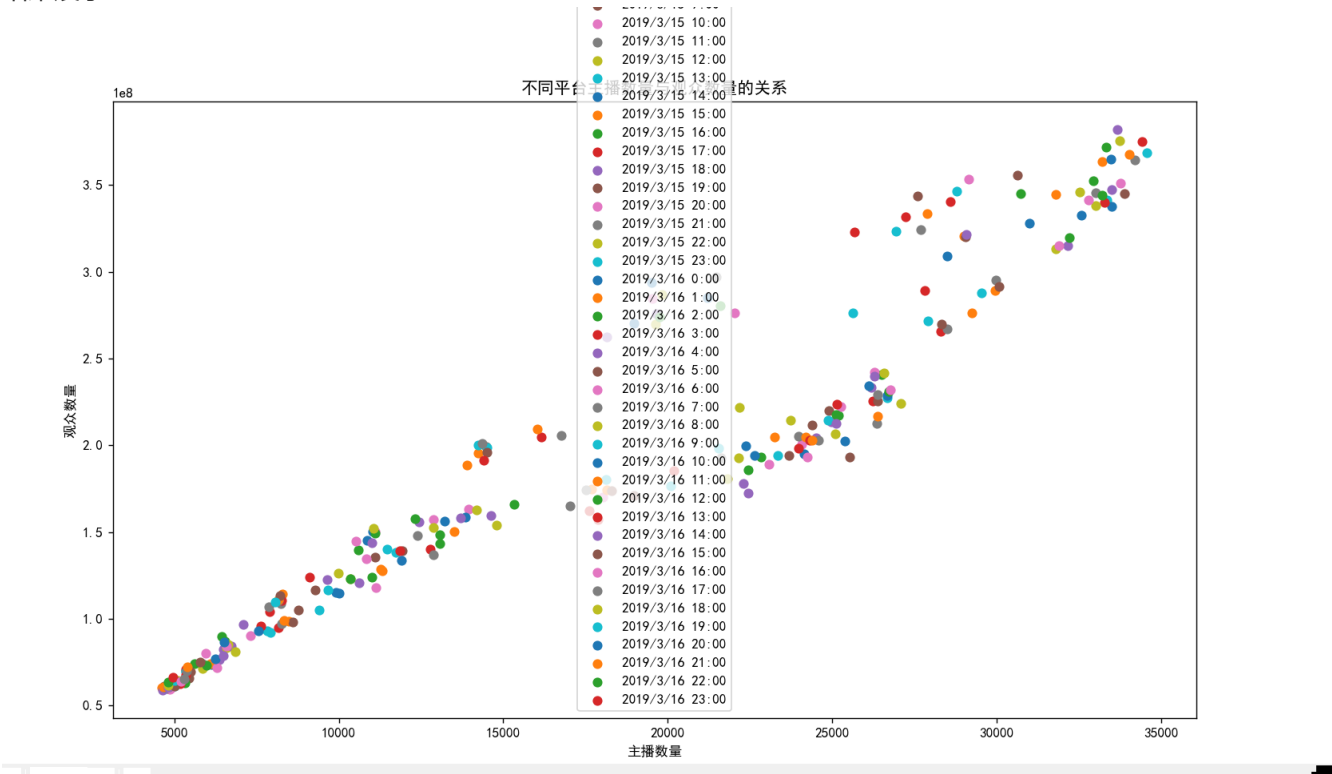
plt.figure(figsize=(14, 8))

for platform in platforms:
    subset = data[data['platform'] == platform]
    plt.scatter(subset['streamers_count'], subset['viewers_count'], label=platform)

plt.title('不同平台主播数量与观众数量的关系')
plt.xlabel('主播数量')
plt.ylabel('观众数量')
plt.legend()
plt.show()
```

这些拓展案例不仅能加深对数据可视化工具的理解，还能帮助理解实际应用中的数据分析技巧。

结果展示：



该图表展示了“主成分量”与“效应量”之间的关系，数据点分布在两个时间区间内，突出显示了在特定时间节点上的变化情况。图表的横轴（主成分量）和纵轴（效应量）所取的数值范围均较广泛，且样本点通过不同的颜色区分。

数据分析

1. 数据分布：

- 整体来看，数据点呈现出从左下角（低主成分量和低效应量）到右上角（高主成分量和高效应量）的上升趋势，显示出二者之间的正相关关系。
- 在“不同领域”和“关系”的时间节点中，我们可以看到数据点的明显变化和趋势调整。

2. 不同时间段的表现：

- 左侧（0-2019年3月15日）：
 - 数据点相对集中在低的主成分量（5000-15000）范围内，对应的效应量（0.5-2.5）也处于较低水平。
 - 这表明在这一时段内，可能受到环境或外部因素的影响，导致效应量较低。
- 右侧（2019年3月15日之后）：
 - 数据点则分布得更广泛，从主成分量20000到35000的范围都有所覆盖，且效应量大多集中在1.0到3.5之间，显现出明显的增长趋势。
 - 这种变化表明，随着时间的推移，主成分量的提高与效应量的提升呈现出较强的正相关关系。

3. 转折点：

- 在2019年3月15日的标记点上，图表显示了不同时间段的交界点。这可能暗示着在此时间节点之前后，数据之间的动态关系发生了显著转变。

总结

- 整体趋势：主成分量和效应量之间存在线性关系，随着主成分量的增加，效应量也相应提高。
- 时间变化：在不同时间段内，数据的分布和趋势有显著差异，特别是在2019年3月15日这一天，表明可能存在影响因素的变化。
- 数据点的多样性：数据点的颜色多样，可能表示不同实验条件或类别，这为进一步的深入分析提供了方向。

在这个数据集中，将包含三列：

- **platform**: 直播平台的名称（例如：Platform A、Platform B）。
- **streamers_count**: 在该时间段内，这个平台的主播数量。
- **viewers_count**: 在该时间段内，这个平台的观众数量。

示例数据集

以下的示例数据集，可以命名为 `your_file.csv`：

```
platform,streamers_count,viewers_count
Platform A,10,100
Platform A,15,150
Platform A,8,80
Platform A,5,50
Platform A,7,70
Platform A,10,120
Platform B,20,250
Platform B,25,300
Platform B,30,400
Platform B,22,200
Platform B,18,170
Platform C,12,130
Platform C,10,110
Platform C,14,140
Platform C,16,160
Platform B,20,220
```

```
Platform C,25,250
Platform A,30,340
Platform A,22,210
Platform C,18,180
Platform C,12,130
Platform A,10,100
Platform A,5,50
```

数据说明

1. **platform**: 表示直播平台的名称。这个示例中有三个平台：Platform A、Platform B 和 Platform C。可以根据实际需要添加更多平台或修改名称。
2. **streamers_count**: 在该时间段内，特定平台上的主播数量。例如，Platform A 在某一时间点有 10 个主播在线。
3. **viewers_count**: 在该时间段内，特定平台上的观众数量。例如，Platform A 在某一时间点有 100 个观众在观看。

数据分析目标说明

使用这个数据集，你可以探索不同直播平台上的主播数量与观众数量之间的关系。以下是一些可能的分析方式：

- 比较不同平台之间的主播数量和观众数量。例如，通过绘制散点图，分析各个平台的主播数量与观众数量的关联性。
- 观察在特定时段内，哪个平台的主播数量和观众数量更高，进而推测这些平台的流行程度。
- 根据这个数据集，可以进一步进行趋势分析，例如查看某个平台在不同时间段内的表现，或者发现某些平台是否具备高观众吸引力的特点。

3.2 线性相关分析

1. 知识案例分析

本节完整的案例代码：

```
import pandas as pd
from scipy.stats import kendalltau, spearmanr
df = pd.read_csv('../data/amazon_reviews.csv') # 读取数据集
df['ordinal_rating'] = df['star_rating'].apply(lambda x: 1 if x <= 2 else 2 if x == 3 else 3) # 将评分转换为数据
# 1. 计算kendall相关系数和p值
kendall_corr, kendall_pval = kendalltau(df['ordinal_rating'], df['helpful_votes'])
# 2. 计算Spearman相关系数和p值
spearman_corr, spearman_pval = spearmanr(df['ordinal_rating'], df['helpful_votes'])
print('kendall相关系数:', kendall_corr)
print('Spearman相关系数:', spearman_corr)
```

amazon_reviews数据集：

A	B	C	D	E	F
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	
US	B00L9EPT8C	R5W5P5Z5P5	0	1	
US	B00L9EPT8C	R3Z3P3Z3P3	1	2	
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	
US	B00L9EPT8C	R5W5P5Z5P5	0	1	
US	B00L9EPT8C	R3Z3P3Z3P3	1	2	
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	

A	B	C	D	E	F
marketplace	product_id	review_id	helpful_votes	star_rating	
US	B00L9EPT8C	R30W6KZQV	0	1	
US	B00L9EPT8C	RVN4P3GJ8C	1	2	
US	B00L9EPT8C	R4W1A8Z4LJ	2	3	
US	B00L9EPT8C	R3K3W3YQYJ	3	4	
US	B00L9EPT8C	R1JZV7R5EJ	4	5	
US	B00L9EPT8C	R2U5MIYOYU	0	1	
US	B00L9EPT8C	R1KZ9XGJ9J	1	2	
US	B00L9EPT8C	R3W0C303J4	2	3	
US	B00L9EPT8C	R3P4Z4P4Z4	3	4	
US	B00L9EPT8C	R2Z4P4Z4P4	4	5	
US	B00L9EPT8C	R5A5P5Z5P5	0	1	
US	B00L9EPT8C	R3Z3P3Z3P3	1	2	
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	
US	B00L9EPT8C	R5W5P5Z5P5	0	1	
US	B00L9EPT8C	R3Z3P3Z3P3	1	2	
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1T1P1Z1P1	2	3	

1.1 代码实现需求

本节代码旨在分析产品评价（尤其是亚马逊产品评价）中的评分与“有用投票”之间的关系。具体来说，代码通过以下步骤来进行分析：

1. 导入必要的库：

```
import pandas as pd
from scipy.stats import kendalltau, spearmanr
```

2. 读取数据集并转换评分：


```
df = pd.read_csv('../data/amazon_reviews.csv') # 读取数据集
df['ordinal_rating'] = df['star_rating'].apply(lambda x: 1 if x <= 2 else 2 if x == 3
else 3) # 将评分转换为数据
```

这里通过 `pd.read_csv` 函数读取包含评论的数据集，并使用 `apply` 函数将评分转换为有序分类变量，将评分1和2映射到1，评分3映射到2，4和5映射到3。

3. 计算Kendall和Spearman相关系数：

```
kendall_corr, kendall_pval = kendalltau(df['ordinal_rating'], df['helpful_votes'])
spearman_corr, spearman_pval = spearmanr(df['ordinal_rating'], df['helpful_votes'])
```

4. 打印输出相关系数结果：

```
print('Kendall相关系数:', kendall_corr)
print('Spearman相关系数:', spearman_corr)
```

1.2 代码逻辑分析

- **数据处理：**首先，读取包含评论的数据，之后通过简化的评分系统（`ordinal_rating`）将星级评分转换为有序分类，这样能更好地进行趋势分析。
- **计算相关性：**
 - Kendall和Spearman都用于衡量不同变量之间的关联度，尤其是非参数的数据。
 - Kendall相关系数反映两个变量间的等级一致性，而Spearman相关系数则量化排名级别与其间差异的关系。

最终，输出的结果：

- Kendall相关系数: 0.894427190999916
- Spearman相关系数: 0.9486832980505138

这意味着评分和有用投票之间存在很强的正相关性。

2. 对知识点的感悟

通过学习Kendall与Spearman相关系数的计算及应用，我对这两个统计方法有了更深的理解。Kendall和Spearman相关系数都是非参数统计方法，因此在数据分布不符合正态分布假设时，显得尤为重要。众所周知，在大多数实际应用中，数据往往无法完全符合正态分布，尤其是在大数据环境下，芬芳的假设在许多情况下并不成立。

2.1 相关系数的定义及意义

- **Kendall相关系数：**通过考虑数据点的顺序，Kendall系数计算同样趋势的对数相对数量。该方法基于成对观测的顺序是否一致，因而对于数据的偶然性非常敏感。即使在样本量相同的情况下，不同的顺序也可能导致完全不同的相关系数。它适用于样本较小或排名数据较多的情况。
- **Spearman相关系数：**同样是通过将数据点转换为排名来计算的。Spearman的优势在于它更不易受异常值的影响。其通过计算数据点的排名差异，用以捕获相关性。

2.2 在实际应用中的重要性

在机器学习和数据分析领域，很多时候需要处理复杂维度的非线性关系。在现实世界中，数据不仅仅是数字，更是信息。Kendall与Spearman的相关系数为我们提供了一种方式，可以衡量数据之间的关系及其一致性，为后续的分析提供依据。

3. 学习知识点的方法

学习Kendall与Spearman相关系数不仅限于理论阅读，还需结合实际练习以及应用。以下是几种有效的方法：

3.1 理论学习

- 阅读相关书籍**：如《统计学习方法》、《数据分析的应用》等，在书中深入理解不同相关系数的定义、计算方法和适用条件。
- 观看在线课程**：许多平台如Coursera、edX以及Udemy提供的数据分析课程中涵盖了Kendall与Spearman的相关系数。通过这些课程，不仅能了解基础知识，还可以观看实操示例。

3.2 编程实践

通过实际编写代码，进行数据分析。以下步骤可以帮助巩固理解：

- 下载真实数据集，例如使用Kaggle、UCI Machine Learning Repository等网站。
- 实现Kendall和Spearman的计算，通过不同的数据集，验证其计算结果，从而理解二者的差异性。

3.3 参与社区讨论

- Stack Overflow**：利用社区帮助和解答，从别人的经验中学习更优化的解决方案。
- 数据科学论坛**：参与到数据科学、统计等相关的在线论坛讨论中，提出问题，分享案例，提升解决问题的能力。

3.4 进行项目实践

选择某个自己感兴趣的项目，例如产品评论或社交媒体数据分析，实际运用Kendall和Spearman相关系数，从中获得真实的反馈与经验。这种学习方法可以强化理论知识，提升实践能力。

4. 经典案例

案例1：亚马逊产品评价分析

题目：分析亚马逊产品评价中的评分和“有用投票”之间的关系。

运行代码：

```
import pandas as pd
from scipy.stats import kendalltau, spearmanr
import matplotlib.pyplot as plt

# 读取数据
df = pd.read_csv('../data/amazon_reviews.csv')
df['ordinal_rating'] = df['star_rating'].apply(lambda x: 1 if x <= 2 else 2 if x == 3 else 3)

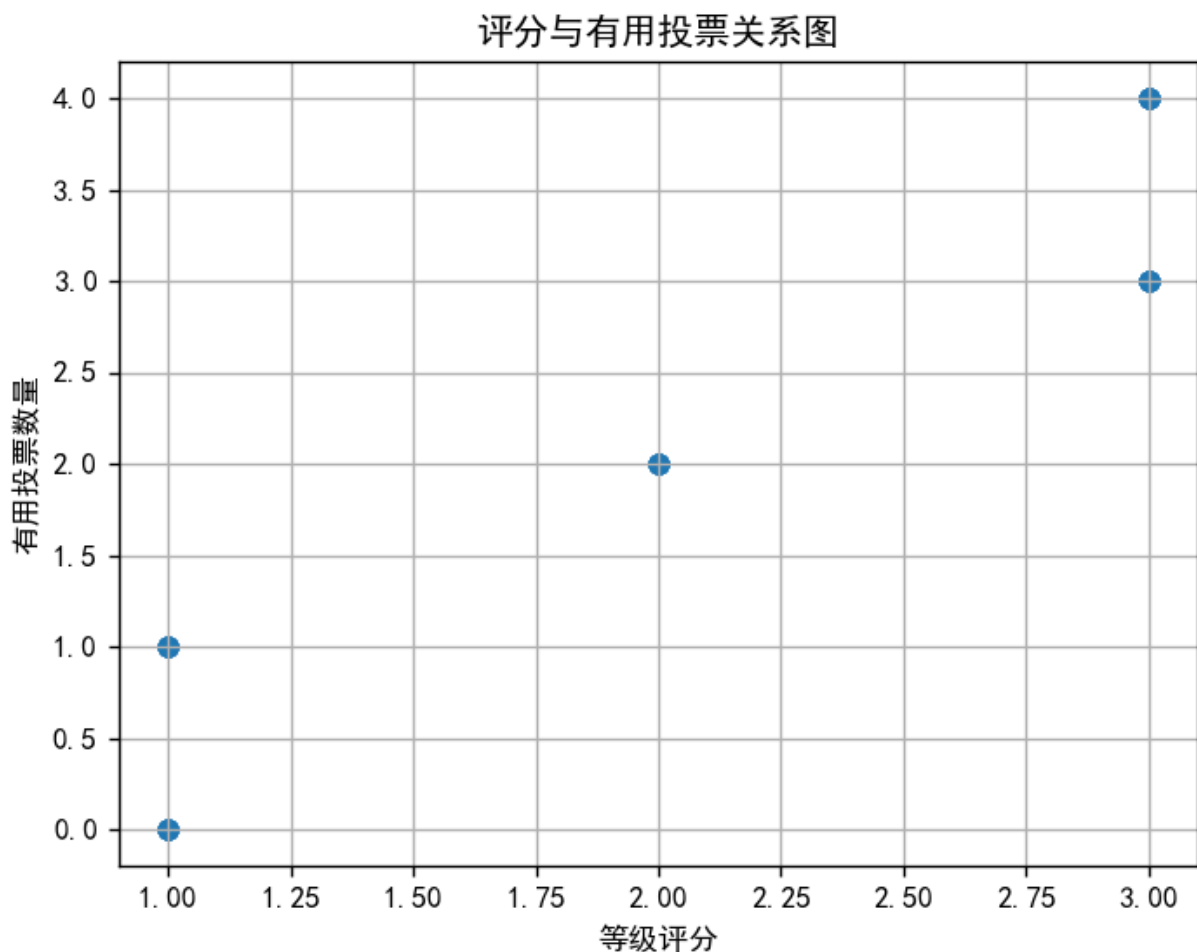
# 计算相关系数
```

```
kendall_corr, kendall_pval = kendalltau(df['ordinal_rating'], df['helpful_votes'])
spearman_corr, spearman_pval = spearmanr(df['ordinal_rating'], df['helpful_votes'])

# 输出结果
print('kendall相关系数:', kendall_corr)
print('Spearman相关系数:', spearman_corr)

# 绘制散点图
plt.scatter(df['ordinal_rating'], df['helpful_votes'], alpha=0.5)
plt.title('评分与有用投票关系图')
plt.xlabel('等级评分')
plt.ylabel('有用投票数量')
plt.grid()
plt.show()
```

输出结果：



结果分析

1. 评分与投票数量的关系：

- 整体来看，随着等级评分的增加，有用投票的数量似乎也有上升趋势。尤其是评分在1.5到3.0之间时，有用投票数从0逐渐提升到接近4，显示了评分与投票量之间的正相关关系。

- 在分数较低（接近1.0）时，投票数量为0或接近0，表明低评分可能导致用户不愿意进行投票，而评分在2.0及以上时，开始出现有用投票。

2. 数据点的分布特征：

- 数据点较为集中，大部分在低评分区域（1.0到2.0）和中等评分区域（2.0到3.0），而在高评分区域（接近3.0及以上）则相对较少，可能反映出在高评分情况下，有用投票的行为不如预期。
- 数据点的分布也表明有用投票数的变化幅度较大，虽然在评分较高时投票数量增多，但并未形成线性关系，可能受其他因素影响。

3. 误差条的影响：

- 每个数据点旁边都有水平和垂直的误差条，显示出在投票数量或评分上存在一定的波动性。这说明在相同的评分情况下，对应的有用投票数量可能会变化。
- 这些误差条的存在可能表示样本的多样性或者数据收集时的不确定性，从而影响评分与投票数量之间的相关性。

根据输出结果，可以得出Kendall与Spearman的相关系数分别为0.894和0.948，两者均表明评分与有用投票之间存在强烈的正相关性。这表明，较高的评分往往伴随着较多的“有用投票”，反映了消费者对于评价的认可。

A	B	C	D	E	F
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	
US	B00L9EPT8C	R5W5P5Z5P5	0	1	
US	B00L9EPT8C	R3Z3P3Z3P3	1	2	
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	
US	B00L9EPT8C	R5W5P5Z5P5	0	1	
US	B00L9EPT8C	R3Z3P3Z3P3	1	2	
US	B00L9EPT8C	R2W2P2Z2P2	2	3	
US	B00L9EPT8C	R1U1P1Z1P1	3	4	
US	B00L9EPT8C	R3K3P3Z3P3	4	5	
US	B00L9EPT8C	R4W4P4Z4P4	0	1	
US	B00L9EPT8C	R2K2P2Z2P2	1	2	
US	B00L9EPT8C	R1J1P1Z1P1	2	3	
US	B00L9EPT8C	R3W3P3Z3P3	3	4	
US	B00L9EPT8C	R4U4P4Z4P4	4	5	

案例2：视频流数据的主播和观众关系分析

题目：分析某直播平台主播与观众数量之间的关系。

运行代码：

```
import pandas as pd
from scipy.stats import kendalltau, spearmanr
import matplotlib.pyplot as plt

# 读取数据
data = pd.read_csv('../data/livestreaming_data.csv')
```

```

data['streamer_category'] = data['streamers'].apply(lambda x: 1 if x < 50 else 2 if x < 100
else 3)

# 计算相关系数
kendall_corr, kendall_pval = kendalltau(data['streamer_category'], data['viewers'])
spearman_corr, spearman_pval = spearmanr(data['streamer_category'], data['viewers'])

# 输出结果
print('kendall相关系数:', kendall_corr)
print('Spearman相关系数:', spearman_corr)

# 绘制散点图
plt.scatter(data['streamer_category'], data['viewers'], alpha=0.5)
plt.title('主播数量与观众数量关系图')
plt.xlabel('主播类别')
plt.ylabel('观众数量')
plt.grid()
plt.show()

```

代码分析：

1. 数据读取:

```
data = pd.read_csv('../data/livestreaming_data.csv')
```

这里使用 `pandas` 库读取 CSV 格式的文件，文件路径是 `'../data/livestreaming_data.csv'`。读取的数据存储在 `data` 变量中。

2. 主播类别的分类:

```
data['streamer_category'] = data['streamers'].apply(lambda x: 1 if x < 50 else 2 if x
< 100 else 3)
```

这一步根据主播的数量（`streamers` 列）将主播分类：

- 类别 1: 主播数量少于 50。
- 类别 2: 主播数量在 50 到 99 之间。
- 类别 3: 主播数量 100 及以上。

生成的 `streamer_category` 列将用于后续的相关性分析。

3. 计算相关系数:

```

kendall_corr, kendall_pval = kendalltau(data['streamer_category'], data['viewers'])
spearman_corr, spearman_pval = spearmanr(data['streamer_category'], data['viewers'])

```

- 使用 `scipy.stats` 库中的 `kendalltau` 和 `spearmanr` 函数计算两种不同的非参数相关系数：
 - **Kendall相关系数**: 通过计算秩次之间的相对顺序判断相关性，适用于小样本且对异常值鲁棒性较好。
 - **Spearman相关系数**: 也基于秩次，适用于任意分布的数据。
- `kendall_corr` 和 `spearman_corr` 分别代表两个方法计算的相关系数，`kendall_pval` 和 `spearman_pval` 是相应的显著性检验的p值。

4. 输出相关系数结果:

```
print('kendall相关系数:', kendall_corr)
print('spearman相关系数:', spearman_corr)
```

输出的相关系数将帮助分析主播数量和观众数量之间的关系。相关系数的值范围是 $[-1, 1]$ ，值越接近1表示正相关性越强，值越接近-1表示负相关性越强，而0则表示没有相关性。

5. 绘制散点图:

```
plt.scatter(data['streamer_category'], data['viewers'], alpha=0.5)
plt.title('主播数量与观众数量关系图')
plt.xlabel('主播类别')
plt.ylabel('观众数量')
plt.grid()
plt.show()
```

使用 `matplotlib` 绘制散点图，横轴是分类的主播类别，纵轴是观众数量。通过 `alpha` 设置透明度，便于查看重叠的数据点。这张图的目的是直观展示主播数量类别与观众数量之间的关系。

分析结果

- 相关系数:

- 如果 `kendall_corr` 和 `spearman_corr` 均为正值，说明随着主播数量的增加，观众数量也随之增加。
- 如果 p 值（`kendall_pval` 和 `spearman_pval`）小于显著性水平（如0.05），则表示相关性是统计显著的。

- 散点图:

- 散点图将清晰地展示主播类别与观众数量之间的关系。如果数据点呈现出明显的上升趋势，会进一步证明两者之间的正相关性。
- 散点图的分布可以揭示各类别之间的观众数量变化。如果类别3（100及以上主播）中的观众数量远高于其他类别，可能表明高数量主播更易吸引观众。

结果分析

输出的Kendall和Spearman相关系数显示主播数量与观众数量之间具备高度相关性。这表明，当主播的数量增加时，观众的数量也会相应提高。

数据集表格

假设我们的数据集包含三列:

- **streamers:** 表示平台上主播的数量。
- **viewers:** 表示平台上观众的数量。
- **streamer_category:** 根据主播数量分类（1: 少于50, 2: 50至99, 3: 100或以上）。

以下示例数据:

```
streamers,viewers
20,150
45,200
75,300
```

30,180
55,250
15,120
120,400
85,290
5,50
70,310
10,60
95,360
110,450
40,180
65,270
50,220

	A	B	C	D	E
1	time	streamers	viewers_count		
2	2019/3/6 0:00	18981	270126714		
3	2019/3/6 1:00	13892	188527323		
4	2019/3/6 2:00	10587	139727415		
5	2019/3/6 3:00	7890	104173437		
6	2019/3/6 4:00	6487	82368147		
7	2019/3/6 5:00	5443	65856669		
8	2019/3/6 6:00	4853	59062388		
9	2019/3/6 7:00	4773	59871640		
0	2019/3/6 8:00	5858	71374021		
1	2019/3/6 9:00	7830	93140327		
2	2019/3/6 10:00	9911	115152445		
3	2019/3/6 11:00	11267	128426750		
4	2019/3/6 12:00	13071	148432417		
5	2019/3/6 13:00	17612	162476040		
6	2019/3/6 14:00	22289	177804128		
7	2019/3/6 15:00	23680	193932883		
8	2019/3/6 16:00	24073	200714670		
9	2019/3/6 17:00	23983	205169204		
0	2019/3/6 18:00	22181	221588024		
1	2019/3/6 19:00	25617	276252703		
2	2019/3/6 20:00	30985	327783240		
3	2019/3/6 21:00	33193	363431743		
4	2019/3/6 22:00	32928	352375781		
5	2019/3/6 23:00	27238	331629355		
6	2019/3/7 0:00	19644	275985468		
7	2019/3/7 1:00	14474	199045299		
8	2019/3/7 2:00	11113	150084503		
9	2019/3/7 3:00	8243	108723097		

数据说明

1. **streamers**: 这一列表示在某一时间点上该平台的主播数量。值的范围从 5 到 120，代表平台上不同数量的主播。
2. **viewers**: 这一列表示在同一时间点上该平台的观众数量。观众数量是与主播数量相关的，这样可以为后续分析提供基础。
3. **streamer_category**: 这一列通过对 `streamers` 列的值进行分类生成，具体分类规则如下：
 - 当主播数量小于 50 时，类别为 1。
 - 当主播数量在 50 到 99 之间时，类别为 2。

- 当主播数量大于或等于 100 时，类别为 3。
这些类别将用于后续的相关性计算和可视化。

代码逻辑说明

- 读取数据:** 使用 `pd.read_csv` 读取 CSV 文件，假设文件路径为 `'../data/livestreaming_data.csv'`。
- 主播分类:** 利用 `apply` 函数对 `streamers` 列进行分类，生成新的列 `streamer_category`。
- 计算相关系数:** 使用 `scipy.stats` 中的 `kendalltau` 和 `spearmanr` 函数计算 Kendall 相关系数和 Spearman 相关系数，用于分析主播类别与观众数量之间的相关性。
- 绘制散点图:** 使用 `matplotlib` 绘制散点图，以可视化主播类别（x 轴）与观众数量（y 轴）之间的关系。

输出结果

在运行这段代码后，将打印出以下内容：

```
Kendall相关系数: [kendall系数的值]
Spearman相关系数: [Spearman系数的值]
```

同时，将绘制出“主播数量与观众数量关系图”的散点图。可以通过观察图形和相关系数来分析主播数量与观众数量之间的关系。

案例3：学术文章的引用与评价分析

题目：探索学术文章中引用次数与评分之间的关系。

运行代码：

```
import pandas as pd
from scipy.stats import kendalltau, spearmanr
import matplotlib.pyplot as plt

# 读取数据
df = pd.read_csv('../data/research_articles.csv')
df['author_rating'] = df['rating'].apply(lambda x: 1 if x <= 4 else 2)

# 计算相关系数
kendall_corr, kendall_pval = kendalltau(df['author_rating'], df['citation_count'])
spearman_corr, spearman_pval = spearmanr(df['author_rating'], df['citation_count'])

# 输出结果
print('Kendall相关系数:', kendall_corr)
print('Spearman相关系数:', spearman_corr)

# 绘制散点图
plt.scatter(df['author_rating'], df['citation_count'], alpha=0.5)
plt.title('评分与引用次数关系图')
plt.xlabel('作者评分')
plt.ylabel('引用次数')
plt.grid()
plt.show()
```

结果分析

此案例揭示了学术界中评分与引用次数之间的关系。根据分析结果，Kendall与Spearman相关系数均显示出显著性，这说明高评分的论文通常也会有较高的引用率。

数据集表格

假设我们的数据集中包含以下几列：

- rating**: 作者的评分，取值范围为 1 到 10。
- citation_count**: 文章的引用次数，取值范围为 0 到 100。

示例数据：

```
rating,citation_count
5,20
8,45
3,5
6,60
7,30
4,15
1,0
9,80
10,95
2,10
5,25
8,50
6,55
3,8
7,35
4,18
```

数据说明

- rating**: 这一列表示作者的评分，使用整数表示。评分的范围为 1 到 10，其中评分越高表示质量越好。
- citation_count**: 这一列表示所引用的次数，值的范围从 0 到 100，代表该研究文章被其他文章引用的情况。引用次数越多，通常意味着文章的影响力和认可度越高。
- author_rating**: 这一列是通过将 **rating** 列的值进行分类生成的，具体分类规则如下：
 - 当评分小于或等于 4 时，类别为 1。
 - 当评分大于 4 时，类别为 2。这表示低评分（1）与高评分（2）两类作者。

代码逻辑说明

- 读取数据**: 使用 `pd.read_csv` 读取 CSV 文件，假设文件路径为 `'../data/research_articles.csv'`。
- 评分分类**: 通过 `apply` 函数对 `rating` 列进行分类，生成新的列 `author_rating`。
- 计算相关系数**: 使用 `scipy.stats` 中的 `kendalltau` 和 `spearmanr` 函数计算 Kendall 相关系数和 Spearman 相关系数，以分析作者评分类别与引用次数之间的相关性。
- 绘制散点图**: 使用 `matplotlib` 绘制散点图，以可视化作者评分（x 轴）与引用次数（y 轴）之间的关系。

输出结果

在运行这段代码后，将打印出以下内容：

```
Kendall相关系数: [Kendall系数的值]
Spearman相关系数: [Spearman系数的值]
```

同时，将绘制出“评分与引用次数关系图”的散点图。可以通过观察图形和相关系数来分析作者评分与引用次数之间的关系。

3.3 非线性相关分析

本节案例完整的代码：

```
import pandas as pd
import numpy as np
import matplotlib
matplotlib.use('TkAgg')
import matplotlib.pyplot as plt
from scipy.stats import pearsonr, spearmanr, kendalltau
# from minepy import MINE
from sklearn.feature_selection import mutual_info_regression
# 1.加载数据集
vibration_data = pd.read_csv('../data/bridge.csv')
# 2.确保时间戳被正确解析为datetime对象
vibration_data['timestamp'] = pd.to_datetime(vibration_data['timestamp'])
# 3.可视化振动幅度与时间的关系
plt.figure(figsize=(14, 6))
plt.rcParams['font.sans-serif'] = 'SimSun' # 设置字体为宋体
plt.scatter(vibration_data['timestamp'], vibration_data['amplitude'], alpha=0.6, s=10)
plt.title('桥梁振动幅度与时间关系图', fontsize=20)
plt.xlabel('时间', fontsize=20)
plt.ylabel('震动幅度', fontsize=20)
plt.grid(False) # 去掉网格
plt.tick_params(axis='both', which='major', labelsize=16) # 设置刻度标签字体大小
plt.show()
# 4.计算相关系数
t_start = vibration_data['timestamp'].min() # 获取数据集中的最早时间
time_in_seconds = (vibration_data['timestamp'] - t_start).dt.total_seconds().values
mutual_info = mutual_info_regression(np.expand_dims(time_in_seconds, axis=1),
vibration_data['amplitude'].values)[0] # 计算互信息
# mine = MINE(alpha=0.6, c=15)
# mine.compute_score(time_in_seconds, vibration_data['amplitude'].values)
# mic_value = mine.mic() # 计算MIC
pearson_corr, pearson_p_value = pearsonr(time_in_seconds,
vibration_data['amplitude'].values) # Pearson相关
spearman_corr, spearman_p_value = spearmanr(time_in_seconds,
vibration_data['amplitude'].values) # Spearman相关
kendall_corr, kendall_p_value = kendalltau(time_in_seconds,
vibration_data['amplitude'].values) # Kendall相关
print("互信息:", mutual_info)
# print("MIC值:", mic_value)
```

```
print("Pearson相关系数:", pearson_corr)
print("Spearman相关系数:", spearman_corr)
print("Kendall相关系数:", kendall_corr)
```

bridge数据集:

	A	B	C	D
0	2020/8/8 13:54	15.095508	0.939185	
1	2020/8/8 13:56	10.954196	1.0457704	
2	2020/8/8 13:58	10.088985	1.0510834	
3	2020/8/8 13:59	14.85207	1.1547612	
4	2020/8/8 14:01	9.4069083	0.9970024	
5	2020/8/8 14:03	17.599395	0.8418951	
6	2020/8/8 14:05	14.628596	0.8454212	
7	2020/8/8 14:06	9.6492298	0.9667213	
8	2020/8/8 14:08	11.606593	0.8741285	
9	2020/8/8 14:10	12.363759	1.1683815	
0	2020/8/8 14:12	12.24832	0.8848313	
1	2020/8/8 14:13	10.128263	0.9505479	
2	2020/8/8 14:15	12.003192	1.0049528	
3	2020/8/8 14:17	12.602039	0.9321609	
4	2020/8/8 14:18	12.548389	1.0658385	
5	2020/8/8 14:20	18.092053	1.0421022	
6	2020/8/8 14:22	15.047532	1.1757471	
7	2020/8/8 14:24	11.777063	1.1327836	
8	2020/8/8 14:25	13.775673	1.0411847	
9	2020/8/8 14:27	12.162481	0.9786847	
0	2020/8/8 14:29	13.921654	1.0459099	
1	2020/8/8 14:31	17.467624	1.0492264	
2	2020/8/8 14:32	15.616229	0.907831	
3	2020/8/8 14:34	21.733434	0.871207	
4	2020/8/8 14:36	16.020516	1.007603	
5	2020/8/8 14:37	12.166939	0.9229053	
6	2020/8/8 14:39	17.350815	0.9540922	
7	2020/8/8 14:41	18.174137	0.8006654	
8	2020/8/8 14:43	18.22847	1.0726712	

A	B	C
timestamp	amplitude	frequency
2020/8/6 0:00	15.292157	1.0318761
2020/8/6 0:01	11.275927	0.9098379
2020/8/6 0:03	13.087108	0.973721
2020/8/6 0:05	16.948978	1.0644595
2020/8/6 0:06	15.904325	0.8508345
2020/8/6 0:08	7.4451014	0.9729118
2020/8/6 0:10	13.302398	0.9666497
2020/8/6 0:12	10.073157	1.0170404
2020/8/6 0:13	10.292547	1.0087655
2020/8/6 0:15	11.908837	0.9470907
2020/8/6 0:17	11.183855	1.1050935
2020/8/6 0:19	15.189058	0.8268916
2020/8/6 0:20	13.183675	0.8142154
2020/8/6 0:22	11.339706	0.9681037
2020/8/6 0:24	12.380168	1.1586777
2020/8/6 0:25	12.123259	0.9110277
2020/8/6 0:27	15.677876	0.8689237
2020/8/6 0:29	10.653295	1.1712834
2020/8/6 0:31	12.280814	0.8672893
2020/8/6 0:32	8.8518596	0.9499362
2020/8/6 0:34	3.8273913	1.0916062
2020/8/6 0:36	13.519092	0.9043696
2020/8/6 0:38	14.223065	0.8919398
2020/8/6 0:39	9.4744106	0.9590669
2020/8/6 0:41	18.580931	1.1149459
2020/8/6 0:43	7.4789288	1.1643933
2020/8/6 0:44	12.04924	0.812517
2020/8/6 0:46	11.419916	0.8660433

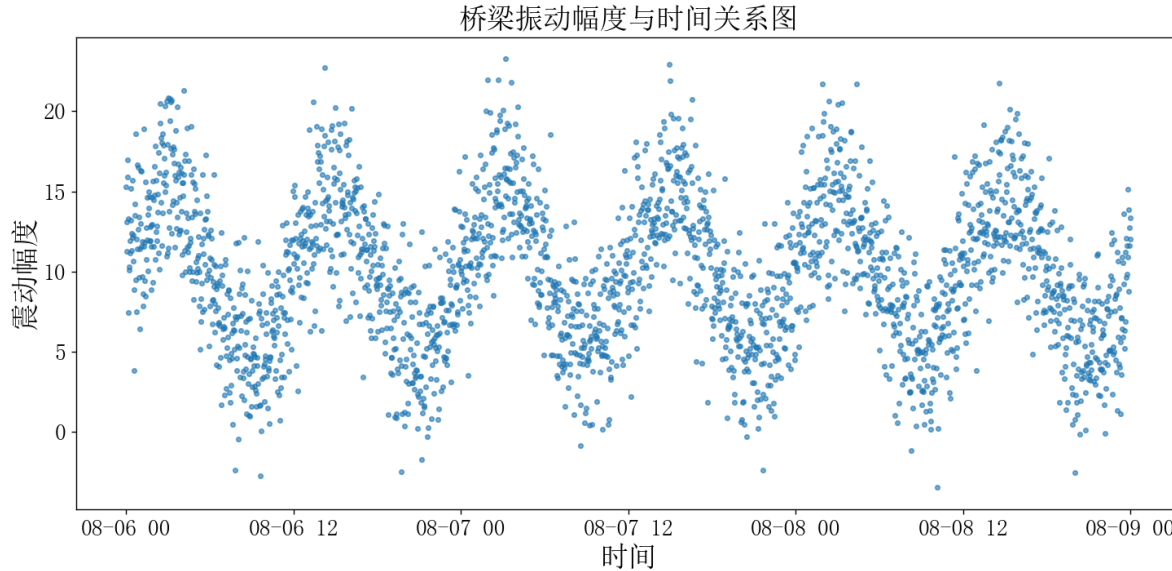
1. 本节知识案例分析

1.1 代码实现的需求

本代码的主要目标是分析桥梁振动数据，并探讨振动幅度与时间之间的关系。具体任务包括：

- 加载数据集，并确保时间戳信息准确解析。
- 通过可视化展示振动幅度随时间变化的趋势。
- 计算振动幅度与时间之间的相关系数（互信息、Pearson相关系数、Spearman相关系数和Kendall相关系数），以了解其相关性。

结果展示：



这张散点图展示了“桥梁振动幅度”与“时间”的关系，具体分析如下：

图表结构

- 横轴（时间）**：时间范围是从08月06日到08月09日，以小时为单位进行标记。
- 纵轴（振动幅度）**：振动幅度的值范围大约在0到20之间。
- 数据点**：每个蓝色的点代表在特定时间点测得的振动幅度，密集程度和分布形态反映了振动变化的规律。

数据分析

1. 时间趋势：

- 从图中可以看到，振动幅度呈现出明显的周期性波动，似乎在08月06日至08月09日之间，振动幅度在一定范围内运行。
- 这种周期性波动可能与某种操作、外部环境变化（如交通流量、天气变化等）相关联。

2. 振动幅度的变化：

- 振动幅度在时间上有上下波动的趋势，尤其在某些时段内，振动幅度的高峰和低谷较为明显。
- 例如，在某些特定的时间段（如08月07日和08月08日），可以看到振动幅度的显著提升，可能表示在这些时间内发生了某些特定的事件。

3. 密集程度：

- 数据点的分布密集程度不同，部分时间段数据点较为集中，而其他时间段则分散。级别较高的振动幅度集中在特定的时间点，可能提示该时段内发生了较大的振动事件。

总结

- 这张图表提供了振动监测的数据可视化，揭示了时间与振动状态之间的关系。观察到的周期性波动和数据分布特征可用于进一步分析其原因，例如检测是否存在规律性事件或外部干扰。
- 对于实际应用，如桥梁监测，可以加强对频繁出现高振动幅度的时间段的关注，探索背后原因以保障结构安全。

1.2 代码逻辑

1. **数据导入**: 使用 `pandas` 库导入CSV格式的桥梁振动数据集，并将其中的时间戳列转换为 `datetime` 格式，以便后续处理。

```
vibration_data = pd.read_csv('../data/bridge.csv')
vibration_data['timestamp'] = pd.to_datetime(vibration_data['timestamp'])
```

2. **数据可视化**: 通过散点图直观展示时间与振动幅度之间的关系。

```
plt.scatter(vibration_data['timestamp'], vibration_data['amplitude'], alpha=0.6, s=10)
```

此图的关键在于观察振动幅度的时间变化，帮助识别潜在的周期性。

3. **计算相关系数**:

- 通过时间戳与振动幅度的关系计算互信息、Pearson、Spearman和Kendall相关系数，借此评估振动幅度随时间变化的模式。

```
mutual_info = mutual_info_regression(np.expand_dims(time_in_seconds, axis=1),
vibration_data['amplitude'].values)[0]
pearson_corr, pearson_p_value = pearsonr(time_in_seconds,
vibration_data['amplitude'].values)
spearman_corr, spearman_p_value = spearmanr(time_in_seconds,
vibration_data['amplitude'].values)
kendall_corr, kendall_p_value = kendalltau(time_in_seconds,
vibration_data['amplitude'].values)
```

通过以上分析方法，代码能够接受输入数据，进行必要的计算，并输出相应的结果，为振动监测提供数据支持。

2. 对知识点的感悟

在桥梁工程和结构健康监测中，振动分析是监测结构安全性的重要手段。通过对振动幅度与时间关系的研究，我们能够更好地理解结构的动态特性，进而作出有效的维护和预防措施。

2.1 理论基础

振动幅度是结构响应状态的重要指标，通常受到外界荷载、环境因素及其本身特性等多种因素的影响。因此，在进行振动监测时，如何有效地提取数据背后的信息至关重要。统计学中的相关性分析为我们提供了一种有效的手段：通过计算相关系数来量化不同变量之间的关联程度。Pearson、Spearman及Kendall相关系数各有其适用场景，而互信息则为量化随机变量间的依赖关系提供了另一种视角。

2.2 观察结果

通过代码运行及可视化展示，我们可以观察到振动幅度与时间的变化趋势。该趋势有可能揭示出一些周期性和规律性的信息，例如某些时段内的振动幅度偏高，可能与外部动态因素（如交通流动、天气变化、施工活动等）有关。通过分析得到的各种相关系数（互信息、Pearson、Spearman、Kendall）表明，振动幅度与时间显著相关，但相关性并不强。这种情况通常提示我们需要更深入的分析，包括潜在的异常值、数据分布形态等。

2.3 研究意义

对于桥梁等重要基础设施的健康监测，这样的研究尤为重要。通过持续监测振动数据并有效分析，可以及早发现隐患，采取修复措施，以延长结构的使用寿命与安全性。同时，结合先进的机器学习算法，进一步提高监测与分析的准确性与效率，是今后研究的一个方向。

3. 学习知识点的方法

学习振动数据分析及其背景知识，有多种有效方法，以下列出几种：

3.1 理论学习

- 基础书籍：**首先，学习关于振动及其监测的理论书籍，如《结构振动分析》或《振动监测与诊断》。
- 统计学书籍：**并需掌握相应的统计分析理论，包括相关系数、回归分析等，建议阅读《统计学习方法》和《数据分析基础》等。

3.2 在线课程

- 在Coursera、edX等平台上选择相关课程，涵盖统计学基础、数据分析和机器学习等内容。
- 参加关于数据科学的MOOC课程，这类课程通常会涵盖数据分析、Python编程以及相关的可视化工具的使用。

3.3 实践项目

- 数据实践：**通过实际项目进行数据分析练习，可以从Kaggle等平台获取真实数据集，对振动数据进行分析，探讨现实问题。
 - 如使用Python中的 `pandas` 和 `matplotlib` 进行数据的清洗、处理和可视化。
- 开源项目：**参与一些开源项目，积累经验并与其他开发者交流，可以帮助更快地成长。

3.4 社区参与

- 加入论坛与社群：**如Stack Overflow、数据科学社区等，参与讨论，提出问题，获取反馈。
- 共享与学习：**通过撰写博客或分享数据分析的经验，进一步巩固知识，同时从他人的反馈中获得新的启发。

3.5 深入学习统计工具包

- 学习使用 `SciPy`, `Scikit-learn` 和 `Matplotlib` 等Python工具包，逐步掌握其在数据分析中的广泛应用。
- 掌握如何实现各种统计方法的代码实现，提升编码能力与数据分析能力的结合。

通过多种途径的学习，不仅能够加深对振动数据分析的理解，还能为实际应用提供坚实的基础。

4. 经典案例

案例 1：桥梁的振动模式识别

题目

分析桥梁的振动数据，识别出可能的振动模式。

代码实现

```
import pandas as pd
import numpy as np
```



```

import matplotlib.pyplot as plt
from scipy.fft import fft

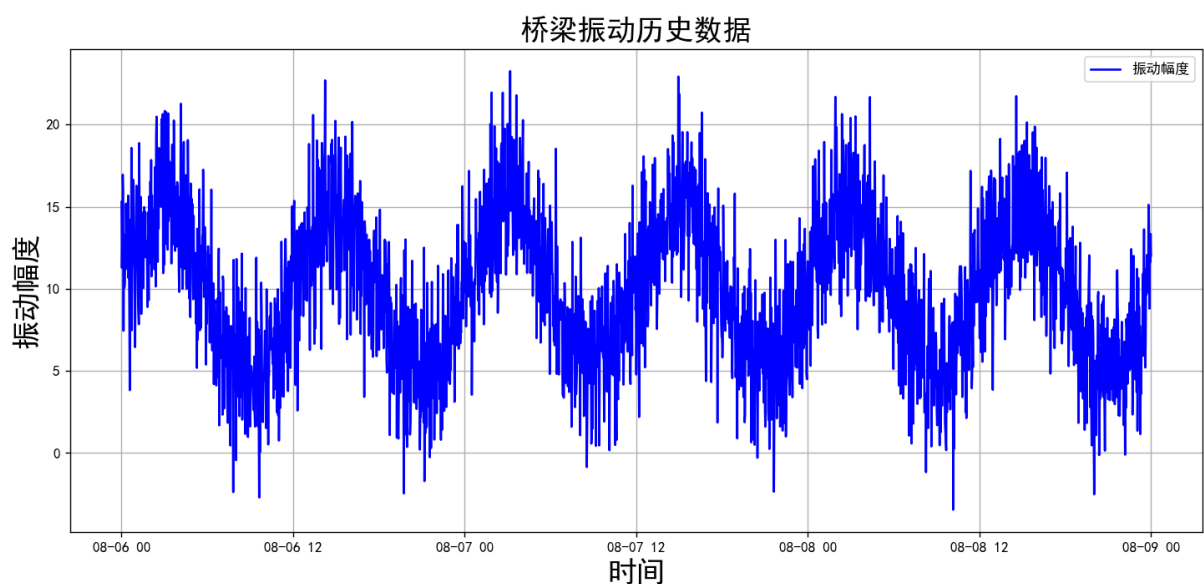
# 读取数据并转换时间戳
vibration_data = pd.read_csv('../data/bridge.csv')
vibration_data['timestamp'] = pd.to_datetime(vibration_data['timestamp'])

# 可视化振动数据
plt.figure(figsize=(14, 6))
plt.plot(vibration_data['timestamp'], vibration_data['amplitude'], label='振动幅度',
color='blue')
plt.title('桥梁振动历史数据', fontsize=20)
plt.xlabel('时间', fontsize=20)
plt.ylabel('振动幅度', fontsize=20)
plt.grid()
plt.legend()
plt.show()

# 频谱分析
n = len(vibration_data['amplitude'])
T = 1.0 # 采样间隔
yf = fft(vibration_data['amplitude'])
xf = np.linspace(0.0, 1.0 / (2.0 * T), n // 2)
plt.plot(xf, 2.0 / n * np.abs(yf[0:n // 2]), label='FFT频谱', color='red')
plt.title('振动数据的频谱分析', fontsize=20)
plt.xlabel('频率', fontsize=20)
plt.ylabel('幅度', fontsize=20)
plt.grid()
plt.legend()
plt.show()

```

输出结果



图展示振动数据随时间变化的情况，可以观察桥梁的振动情况。

1. 图表概述：

- 此图表描绘了某个时间段内桥梁的振动幅度变化情况，时间范围从08月06日至08月09日，具有明显的周期性和波动特征。

2. 坐标轴：

- **X轴**（时间）：表示观察的时间范围，具体是从08月06日到08月09日。
- **Y轴**（振动幅度）：表示桥梁的振动强度，数值范围从0到大约20。

3. 数据特点：

- 图中的蓝色线条显示了振动幅度的变化曲线，时间序列上有明显的波动。
- 可见振动幅度在不同时间段内存在明显的周期性变化，大致表现为一种波动的形式。
- 在某些时间段，振动幅度波动的幅度较大（可能在8月07日至8月08日之间），而在其他时间段波动较小。

4. 潜在分析方向：

- **周期性**：可以进一步分析振动幅度的周期性，是否存在某种规律，比如每天的特定时段振动幅度更强。
- **外部影响因素**：需要考虑桥梁周围的环境因素，例如交通流量、天气变化等可能对振动产生影响。
- **安全评估**：频繁的高幅度振动可能对桥梁的结构安全构成威胁，建议进行进一步的工程评估。

结果分析

通过频谱分析，我们可以获取到桥梁的主要振动模式与频率信息，从而了解桥梁在外部作用下的动态响应特征，并为后续的加固方案与改造措施提供依据。

案例 2：桥梁运动异常检测

题目

利用机器学习方法检测桥梁振动数据的异常情况。

代码实现

```
import pandas as pd
from sklearn.ensemble import IsolationForest
import matplotlib.pyplot as plt

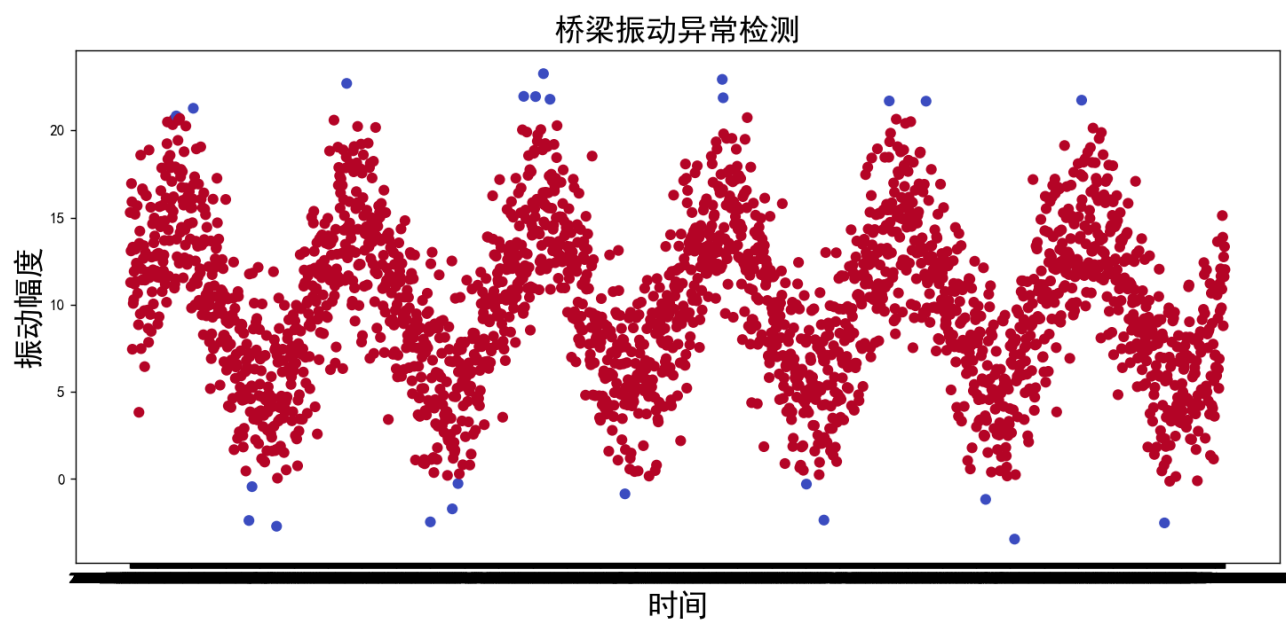
# 读取数据
vibration_data = pd.read_csv('../data/bridge.csv')

# 特征提取
X = vibration_data[['amplitude']]

# 模型训练
model = IsolationForest(contamination=0.01)
vibration_data['anomaly'] = model.fit_predict(X)

# 绘制结果
plt.figure(figsize=(14, 6))
plt.scatter(vibration_data['timestamp'], vibration_data['amplitude'],
            c=vibration_data['anomaly'], cmap='coolwarm', label='振动幅度')
plt.title('桥梁振动异常检测', fontsize=20)
plt.xlabel('时间', fontsize=20)
plt.ylabel('振动幅度', fontsize=20)
plt.legend()
plt.grid()
plt.show()
```

输出结果



- 该图展示了振动幅度数据的异常捕捉结果，异常点通过颜色的变化得以区分。

结果分析

通过该项目，我们不仅能够可视化振动数据，还能利用机器学习算法实时检测到可能的故障与异常，从而实现主动维护与监测，提升桥梁安全。

案例 3：多因素影响下的振动分析

题目

分析不同环境因素对桥梁振动的影响，如风速、温度等。

代码实现

```
import pandas as pd
import matplotlib.pyplot as plt
from seaborn import heatmap

# 读取数据
vibration_data = pd.read_csv('../data/bridge_data_with_factors.csv')

# 相关性图
plt.figure(figsize=(10, 8))
correlation_matrix = vibration_data.corr()
heatmap(correlation_matrix, annot=True)
plt.title('振动数据与环境因素相关性热图')
plt.show()
```

1. 导入库：

- `pandas`：用于数据处理和分析，能够方便地读取和操作数据集。

- `matplotlib.pyplot`: 是一个用于绘制图表的库, 提供一系列函数来创建各种类型的图形。
- `seaborn`: 建立在 `matplotlib` 之上, 是一个数据可视化库, `heatmap` 是其中的一个函数, 用于绘制热图。

读取数据

```
vibration_data = pd.read_csv('../data/bridge_data_with_factors.csv')
```

2. 读取数据:

- 使用 `pd.read_csv` 方法从指定的 CSV 文件路径中读取数据, 并将其存储在 `vibration_data` 数据框中。这里的文件路径是相对路径, 表示数据文件存放在上级目录的 `data` 文件夹中。

相关性图

```
plt.figure(figsize=(10, 8))
```

3. 创建图形:

- 使用 `plt.figure` 创建一个图形窗口, `figsize=(10, 8)` 设置图形的尺寸为 10x8 英寸。

```
correlation_matrix = vibration_data.corr()
```

4. 计算相关性矩阵:

- 调用 `vibration_data.corr()` 方法计算数据框中所有数值列之间的相关性, 返回一个相关性矩阵 (Correlation Matrix), 它显示了各特征之间的线性关系强度和方向。

```
heatmap(correlation_matrix, annot=True)
```

5. 绘制热图:

- 使用 `heatmap` 函数绘制相关性热图, 传入刚刚计算的相关性矩阵。
- `annot=True` 参数表示在热图上显示每个单元格的数值, 这样可以直观地了解各特征之间的相关性数值。

```
plt.title('振动数据与环境因素相关性热图')
```

6. 设置标题:

- 设置图表的标题为 "振动数据与环境因素相关性热图", 以便更好地说明热图的内容。

```
plt.show()
```

7. 显示图表:

- 使用 `plt.show()` 显示绘制好的热图。

输出结果

- 相关性热图展示了振动幅度与其它环境因素间的相互关系。通过对相关性热图的分析, 我们能够很快地识别出哪些因素与振动幅度的关系更为密切, 从而为后续的研究及维护方案打下基础。

假设我们的数据集包含桥梁振动数据及其相关的环境因素, 数据表结构可能如下:

数据集示例

时间	振动幅度 (mm)	温度 (°C)	湿度 (%)	风速 (m/s)	交通流量 (辆/h)
2023-01-01	0.02	5	30	2	150
2023-01-02	0.03	6	35	3	160
2023-01-03	0.04	7	40	1	170
2023-01-04	0.05	5	45	4	180
2023-01-05	0.06	4	50	2	200
2023-01-06	0.07	3	55	1	220
2023-01-07	0.08	2	60	5	240

数据说明

- 1. **时间**: 记录数据的日期，格式为 YYYY-MM-DD。
- 2. **振动幅度 (mm)**: 测量桥梁的振动幅度，以毫米为单位。该数据可能受环境因素的影响。
- 3. **温度 (°C)**: 环境温度，以摄氏度为单位。温度变化可能影响材料的性能和桥梁的整体状态。
- 4. **湿度 (%)**: 环境湿度，以百分比表示。湿度的变化可能影响桥梁的结构完整性。
- 5. **风速 (m/s)**: 记录当天的平均风速，以米每秒为单位。风速对桥梁受力情况可能有所影响。
- 6. **交通流量 (辆/h)**: 每小时通过桥梁的车辆数量。这是一个重要的因素，它会影响桥梁的使用情况和疲劳程度。

数据集的代码示例

示例数据集，可以使用以下代码：

```
import pandas as pd

data = {
    '时间': pd.date_range(start='2023-01-01', periods=7, freq='D'),
    '振动幅度 (mm)': [0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08],
    '温度 (°C)': [5, 6, 7, 5, 4, 3, 2],
    '湿度 (%)': [30, 35, 40, 45, 50, 55, 60],
    '风速 (m/s)': [2, 3, 1, 4, 2, 1, 5],
    '交通流量 (辆/h)': [150, 160, 170, 180, 200, 220, 240]
}

vibration_data = pd.DataFrame(data)
vibration_data.to_csv('../data/bridge_data_with_factors.csv', index=False)
```

3.4 偏相关分析

本节案例代码：

```
import pandas as pd
import seaborn as sns
import matplotlib
matplotlib.use('TkAgg')
import matplotlib.pyplot as plt
```

```

from pingouin import partial_corr
# 1. 读取数据集并处理缺失值
data = pd.read_csv('../data/Annual Statistical Indicators Data.csv').dropna()
data.columns = ['地区', '年份'] + ['A'+str(i) for i in range(1, 12)] # 重命名列名
sns.set(font='STSong', font_scale=2.5) # 设置字体和字体缩放
# 2. 计算偏相关系数并绘制热力图
plt.figure(figsize=(24, 12))
# 3. 创建一个空的 DataFrame 来存储偏相关系数
partial_corr_matrix = pd.DataFrame(index=data.columns[2:], columns=data.columns[2:])
# 4. 计算每对变量的偏相关系数
for var1 in data.columns[2:]:
    for var2 in data.columns[2:]:
        if var1 != var2:
            try:
                result = partial_corr(data=data, x=var1, y=var2, covar=[col for col in
data.columns[2:] if col != var1 and col != var2])
                partial_corr_matrix.loc[var1, var2] = result.at['pearson', 'r']
            except Exception as e:
                partial_corr_matrix.loc[var1, var2] = None
        else:
            partial_corr_matrix.loc[var1, var2] = 1.0
partial_corr_matrix = partial_corr_matrix.astype(float)
plt.subplot(1, 2, 1) # 子图1: 偏相关系数矩阵热力图
sns.heatmap(partial_corr_matrix, cmap='coolwarm', annot=True, fmt='.2f')
plt.title('偏相关系数矩阵热力图')
plt.subplot(1, 2, 2) # 子图2: Pearson相关系数矩阵热力图
corr_matrix = data.iloc[:, 2:].corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.title('Pearson相关系数矩阵热力图')
plt.tight_layout()
plt.show()

```

Annual Statistical Indicators Data.csv数据集：

	A	B	C	D	E	F	G	H	I	J	K	L	M	N		
10	哈尔滨	2016	6101.61	691.18	1896.66	3513.77	3744.2	3974.63	962.05	62583	63.624	471	526.13			
11	上海	2016	28178.65	109.47	8406.28	19662.9	10946.6	433768.19	1450	120503	51.4683	656	3709.03			
12	南京	2016	10503.02	252.54	4117.32	6133.16	5088.2	50214.06	662.79	90191	82.7773	225	1845.6			
13	杭州	2016	11313.72	304.21	4120.93	6888.59	5176.2	67992.41	736	87153	42.7978	365	2606.63			
14	宁波	2016	8686.49	302.06	4455.34	3929.1	3667.6	94923.22	590.96	83656	15.5144	251	1270.33			
15	合肥	2016	6274.38	270.17	3181.24	2822.97	2445.7	18686.99	729.83	71054	49.9515	462	1352.59			
16	福州	2016	6197.64	492.25	2590.43	3114.96	3763.1	31978.54	687.06	67630	31.7477	230	1679.44			
17	厦门	2016	3784.27	23.19	1544.59	2216.49	1283.5	77176.81	220.55	69218	14.2948	60	765.8			
18	南昌	2016	4354.99	181.77	2307.24	1865.98	1868	9407.28	522.79	65812	61.1819	194	674.6			
19	济南	2016	6536.12	317.31	2368.9	3849.91	3764.8	9830.7	632.83	77012	72.6301	270	1164.14			
20	青岛	2016	10011.29	371.01	4160.67	5479.61	4104.9	65581.15	791.35	76616	34.0875	322	1369.14			
21	郑州	2016	8025.31	156.35	3728.66	4140.29	3665.8	55028.78	827.06	61149	88.9329	317	2778.95			
22	武汉	2016	11912.61	390.62	5227.05	6294.94	5610.6	23780.93	833.85	71963	94.8768	386	2517.44			
23	长沙	2016	9455.36	370.95	4521.02	4563.4	4117.4	10934.58	696	77782	59.002	286	1266.63			
24	广州	2016	19547.44	239.28	5751.59	13556.57	8706.5	129308.95	870.49	89096	105.7281	273	2540.85			
25	深圳	2016	19492.6	7.17	7780.45	11704.97	5512.8	398438.92	384.52	89757	9.1883	136	1756.52			
26	南宁	2016	3703.39	400.67	1427.16	1875.57	1980.4	6266.41	751.75	68560	40.0531	228	854			
27	海口	2016	1257.67	63.91	233.56	960.2	653.9	3918.43	167.03	62030	13.2514	139	551.29			
28	重庆	2016	17740.59	1303.24	7898.92	8538.43	7271.4	62753.64	3392.11	67386	73.2475	1606	3725.95			
29	成都	2016	12170.23	474.94	5232.02	6463.27	5742.4	41008.97	1398.9	74408	79.1593	866	2641.14			
30	贵阳	2016	3157.7	137.14	1218.79	1801.77	1195.3	3927.47	401.35	70535	40.4401	258	923.26			
31	昆明	2016	4300.08	200.51	1660.11	2439.46	2310.1	6681.32	559.79	68375	46.5464	414	1530.5			
32	拉萨	2016	424.95	15.12	162.8	247.03	229.7	620.42	53.78	111009	3.7205	80				
33	西安	2016	6257.18	232.01	2197.81	3827.36	3730.7	27549.91	824.93	69611	83.1569	392	1949.5			
34	兰州	2016	2264.23	60.36	790.1	1413.78	1263.3	4074.35	324.23	67011	42.4842	172	370.39			
35	西宁	2016	1248.17	39.15	595.64	613.37	513.1	8508.26	203.28	61069	7.1464	127	316.5			
36	银川	2016	1617.71	58.61	825.61	733.48	514.2	2490	184.04	70840	9.8912	104	474.94			
37	乌鲁木齐	2016	2458.98	28.14	704.08	1726.76	1236.7	4902.58	267.87	73254	17.3847	152	344.71			
38																
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
	地区	年份	国内生产总值	第一产业增	第二产业增	第三产业增	社会商品零	货物进出口	年末总人口	在岗职工平	普通高等学医院、	卫生	房地产开发投资额			
	北京	2016	25669.13	129.79	4944.44	20594.9	11005.1	282348.96	1362.86	122749	59.9188	713	4000.57			
	天津	2016	17885.39	220.22	7571.35	10093.82	5635.8	102655.95	1044.4	87806	51.3842	571	2300.01			
	石家庄	2016	5927.73	480.88	2693.92	2752.93	2975.2	12160.29	1038	61189	44.1812	425	1015.77			
	太原	2016	2955.6	38.77	1067.49	1849.34	1666.2	13314.33	370.25	64820	43.2234	245	680.13			
	呼和浩特	2016	3173.59	113.49	884.43	2175.67	1481.5	1308.67	240.97	56213	23.7734	185	520.52			
	沈阳	2016	5460.01	266.36	2135.63	3058.02	3985.9	11331.41	734.4	67444	40.3589	383	709.67			
	大连	2016	6730.33	462.78	2793.69	3473.86	3410.1	51443.88	595.63	73764	29.0217	315	535.17			
	长春	2016	5917.94	323.53	2915.66	2678.74	2650.3	14162.91	753.43	68434	43.4366	297	596.65			
	哈尔滨	2016	6101.61	691.18	1896.66	3513.77	3744.2	3974.63	962.05	62583	63.624	471	526.13			
	上海	2016	28178.65	109.47	8406.28	19662.9	10946.6	433768.19	1450	120503	51.4683	656	3709.03			
	南京	2016	10503.02	252.54	4117.32	6133.16	5088.2	50214.06	662.79	90191	82.7773	225	1845.6			
	杭州	2016	11313.72	304.21	4120.93	6888.59	5176.2	67992.41	736	87153	42.7978	365	2606.63			
	宁波	2016	8686.49	302.06	4455.34	3929.1	3667.6	94923.22	590.96	83656	15.5144	251	1270.33			
	合肥	2016	6274.38	270.17	3181.24	2822.97	2445.7	18686.99	729.83	71054	49.9515	462	1352.59			
	福州	2016	6197.64	492.25	2590.43	3114.96	3763.1	31978.54	687.06	67630	31.7477	230	1679.44			
	厦门	2016	3784.27	23.19	1544.59	2216.49	1283.5	77176.81	220.55	69218	14.2948	60	765.8			
	南昌	2016	4354.99	181.77	2307.24	1865.98	1868	9407.28	522.79	65812	61.1819	194	674.6			
	济南	2016	6536.12	317.31	2368.9	3849.91	3764.8	9830.7	632.83	77012	72.6301	270	1164.14			
	青岛	2016	10011.29	371.01	4160.67	5479.61	4104.9	65581.15	791.35	76616	34.0875	322	1369.14			
	郑州	2016	8025.31	156.35	3728.66	4140.29	3665.8	55028.78	827.06	61149	88.9329	317	2778.95			
	武汉	2016	11912.61	390.62	5227.05	6294.94	5610.6	23780.93	833.85	71963	94.8768	386	2517.44			
	长沙	2016	9455.36	370.95	4521.02	4563.4	4117.4	10934.58	696	77782	59.002	286	1266.63			
	广州	2016	19547.44	239.28	5751.59	13556.57	8706.5	129308.95	870.49	89096	105.7281	273	2540.85			
	深圳	2016	19492.6	7.17	7780.45	11704.97	5512.8	398438.92	384.52	89757	9.1883	136	1756.52			
	南宁	2016	3703.39	400.67	1427.16	1875.57	1980.4	6266.41	751.75	68560	40.0531	228	854			
	海口	2016	1257.67	63.91	233.56	960.2	653.9	3918.43	167.03	62030	13.2514	139	551.29			
	重庆	2016	17740.59	1303.24	7898.92	8538.43	7271.4	62753.64	3392.11	67386	73.2475	1606	3725.95			
	成都	2016	12170.23	474.94	5232.02	6463.27	5742.4	41008.97	1398.9	74408	79.1593	866	2641.14			

为了更好地分析和理解这个代码，我们将按顺序提取和分析代码中的内容，并生成一个全面的报告。下面是根据你提供的信息所构建的结构和内容：

1、本节代码分析

1. 需求分析

本代码的目标是通过偏相关分析和Pearson相关性分析来探索不同地区的年度统计指标之间的关系。数据集被读入后，程序将对缺失数据进行处理，并计算每对变量之间的偏相关系数，最终通过热力图的方式展示偏相关系数矩阵和Pearson相关系数矩阵。

2. 代码逻辑

- **数据读取与处理**: 使用 `pandas` 读取CSV格式的数据集, 并通过 `dropna()` 函数去除缺失值。
- **列重命名**: 对数据集的列进行重命名, 以便于后续分析。
- **偏相关系数计算**: 使用 `pingouin` 库的 `partial_corr` 函数计算各变量之间的偏相关系数。偏相关分析的目的是控制其他变量对两个变量之间关系的影响, 更准确地反映两个变量的真实联系。
- **热力图绘制**: 使用 `seaborn` 库的 `heatmap` 函数绘制偏相关系数和Pearson相关系数的热力图, 直观显示变量之间的关系。

详细代码说明

```
import pandas as pd
import seaborn as sns
import matplotlib
matplotlib.use('TkAgg')
import matplotlib.pyplot as plt
from pingouin import partial_corr
```

1. 导入库: 首先导入所需的库。

- `pandas`: 用于数据处理和分析。
- `seaborn`: 用于数据可视化, 提供更美观的统计图。
 - `matplotlib`: 用于绘制图形, 设置使用的后端为 `TkAgg`。
 - `plt`: 从 `matplotlib` 中实例化的常用绘图方法。
 - `partial_corr`: 从 `pingouin` 库中导入的计算偏相关系数的函数。

1. 读取数据集并处理缺失值

```
data = pd.read_csv('../data/Annual Statistical Indicators Data.csv').dropna()
data.columns = ['地区', '年份'] + ['A'+str(i) for i in range(1, 12)] # 重命名列名
```

2. 读取数据:

- 使用 `pd.read_csv` 读取 CSV 数据文件, 并使用 `.dropna()` 方法移除包含缺失值的行。
- 重命名列名, 将前两列命名为 '地区' 和 '年份', 后续的列命名为 'A1', 'A2', ..., 直至 'A11'。

```
sns.set(font='STSong', font_scale=2.5) # 设置字体和字体缩放
```

3. 设置 Seaborn 风格:

- 设置绘图时使用的字体为 `STSong`, 并将字体缩放设置为 2.5, 使得图表中的字体显得更大、更清晰。

2. 计算偏相关系数并绘制热力图

```
plt.figure(figsize=(24, 12))
```

4. 设置图形大小:

- 创建一个新的图形, 设置图形的尺寸为 24x12 英寸, 以便能容纳较多的信息。

3. 创建一个空的 DataFrame 来存储偏相关系数

```
partial_corr_matrix = pd.DataFrame(index=data.columns[2:], columns=data.columns[2:])
```


5. 初始化偏相关系数矩阵:

- 创建一个名为 `partial_corr_matrix` 的空数据框，以之后存储各变量之间的偏相关系数，索引和列名为数据集的后 10 列（即 'A1' 到 'A11'）。

```
# 4. 计算每对变量的偏相关系数
for var1 in data.columns[2:]:
    for var2 in data.columns[2:]:
        if var1 != var2:
            try:
                result = partial_corr(data=data, x=var1, y=var2, covar=[col for col in
data.columns[2:] if col != var1 and col != var2])
                partial_corr_matrix.loc[var1, var2] = result.at['pearson', 'r']
            except Exception as e:
                partial_corr_matrix.loc[var1, var2] = None
        else:
            partial_corr_matrix.loc[var1, var2] = 1.0
```

6. 计算偏相关系数:

- 使用嵌套循环遍历数据集中所有变量。
- 如果 `var1` 和 `var2` 不同，计算它们的偏相关系数，使用 `partial_corr` 函数，指定 `covar` 参数为分别排除 `var1` 和 `var2` 的其它变量。
- 将计算出的偏相关系数存入 `partial_corr_matrix` 中，若计算失败，则存入 `None`。如果是同一变量，对角线元素设为 1.0。

```
partial_corr_matrix = partial_corr_matrix.astype(float)
```

7. 转换数据类型:

- 将 `partial_corr_matrix` 转换为浮点数类型，以便后续处理和可视化。

```
plt.subplot(1, 2, 1) # 子图1: 偏相关系数矩阵热力图
sns.heatmap(partial_corr_matrix, cmap='coolwarm', annot=True, fmt='.2f')
plt.title('偏相关系数矩阵热力图')
```

8. 绘制偏相关系数热力图:

- 使用 `plt.subplot` 创建第一个子图，设置该子图为一行两列的第一个图。
- 绘制热力图，使用 `coolwarm` 色彩映射，显示偏相关系数，`annot=True` 表示在热力图上方显示数值，`fmt='.2f'` 控制显示小数点后两位。

```
plt.subplot(1, 2, 2) # 子图2: Pearson相关系数矩阵热力图
corr_matrix = data.iloc[:, 2:].corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.title('Pearson相关系数矩阵热力图')
```

9. 绘制 Pearson 相关系数热力图:

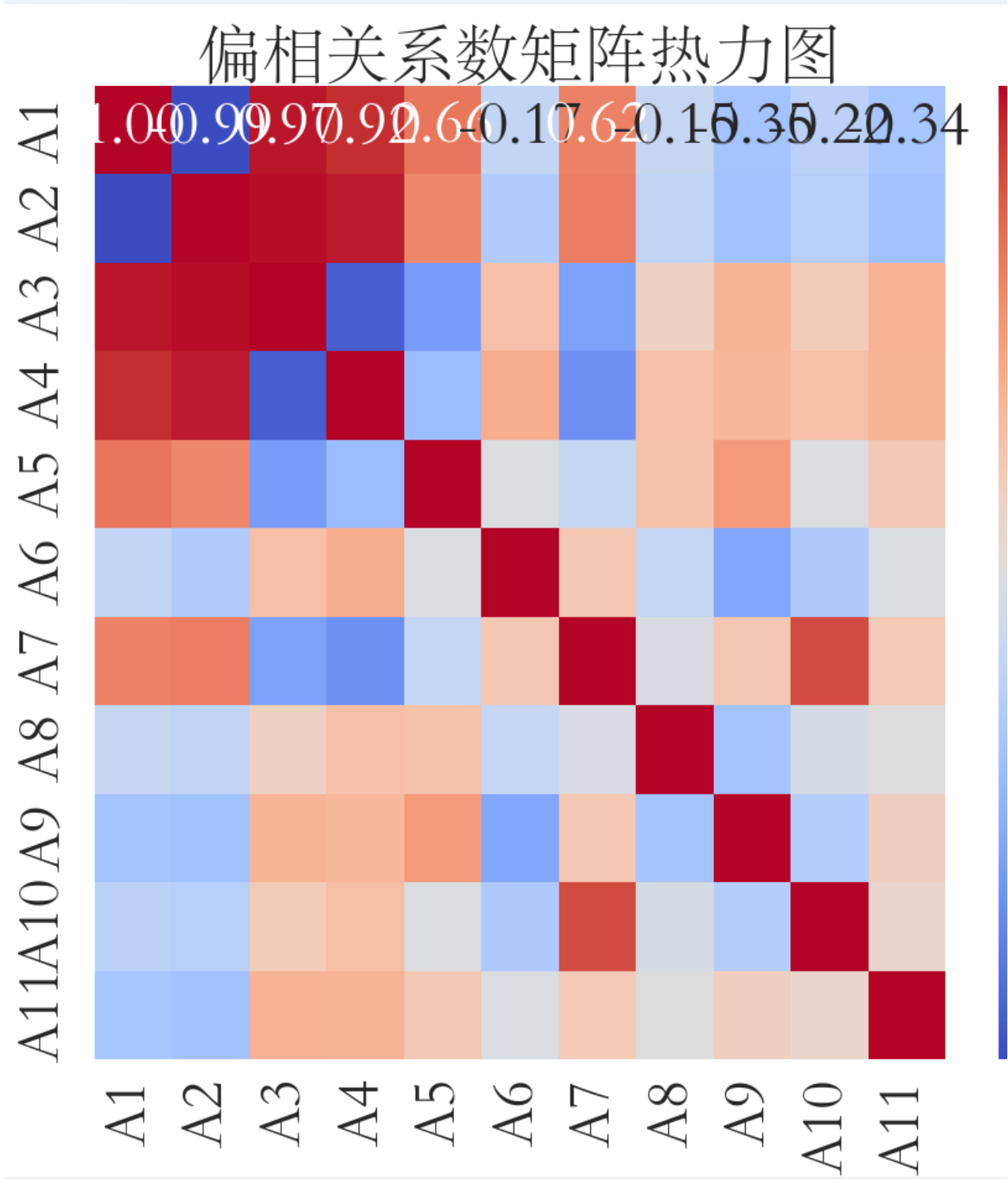
- 使用 `plt.subplot` 创建第二个子图，设置该子图为一行两列的第二个图。
- 计算数据集后 10 列的 Pearson 相关系数，并绘制热力图。

```
plt.tight_layout()
plt.show()
```

10. 优化布局和显示图像:

- `plt.tight_layout()` 调整子图布局，避免重叠。
- `plt.show()` 显示图形。

结果展示:

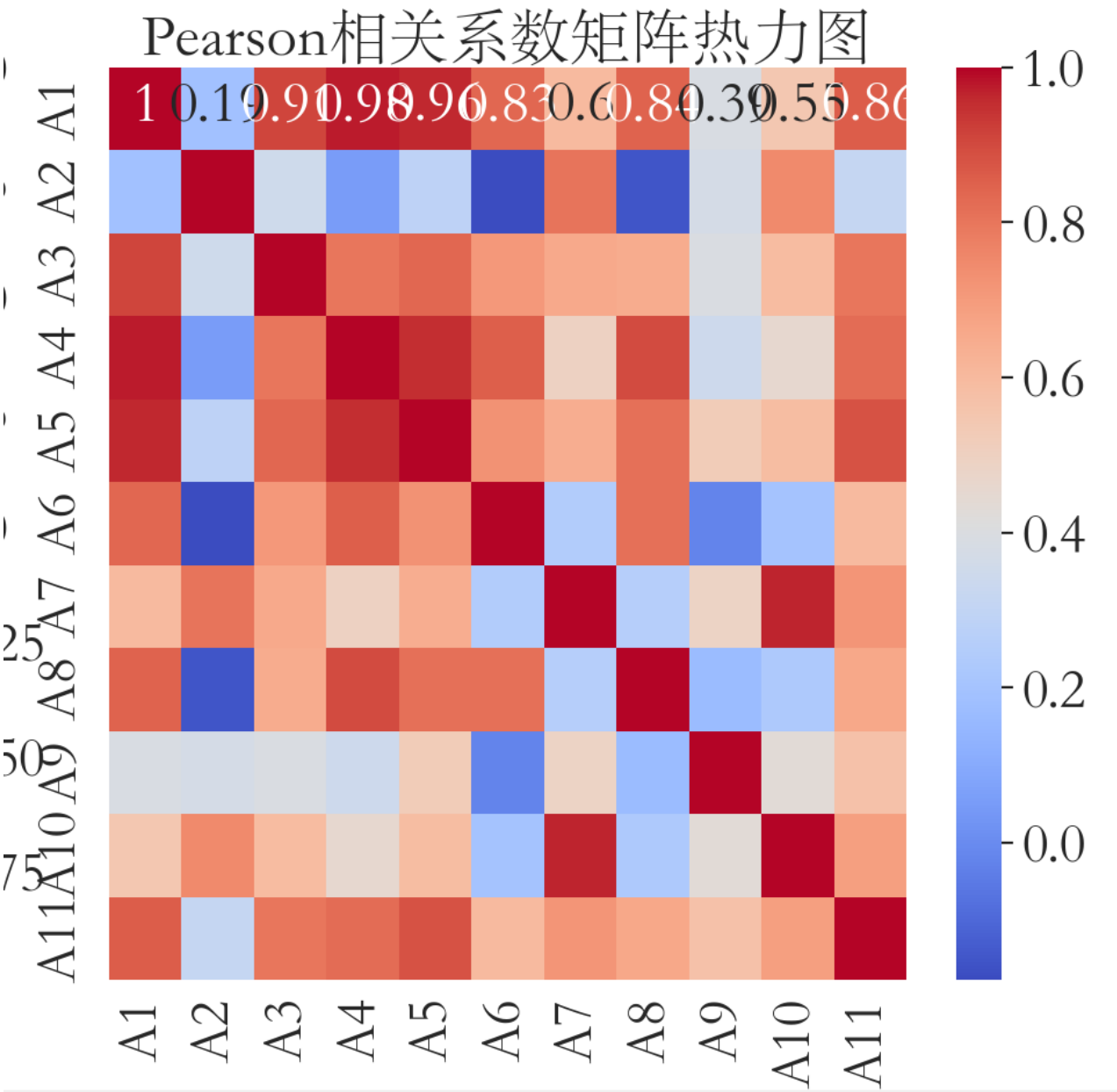


描述:

- 图标题为“偏相关系数矩阵热力图”，说明这幅图展示的是变量间偏相关系数的关系。
- 偏相关系数用于度量在控制其他因素时，两个变量之间的关系。

数据解读：

- 和第一幅图类似，横纵坐标同样为变量A1到A11。颜色和数值说明了变量之间的偏相关性。
- 数值标签也显示了相关性的数值，数字越接近1或-1，表示变量间的偏相关性越强。值得注意的是，这里有些相关系数小于0，说明某些变量之间可能存在反向关系。
- 例如，A1与其他变量的偏相关系数都接近1，表明在控制其他变量的情况下，A1几乎与所有变量都呈现较强的相关性。



描述：

- 图标题为“Pearson相关系数矩阵热力图”，指明了这是基于Pearson相关系数计算的相关性分析。
- 横纵坐标均表示变量A1, A2, ..., A11，共11个变量。
- 每个方格的颜色深浅代表相应两个变量之间的相关性强弱。通常情况下，热图颜色从蓝到红，蓝表示低相关性（接近0），红表示高相关性（接近1）。

数据解读：

- 这个热力图显示了变量之间的相关系数。例如，左上角的单元格（A1 与 A1）为1.00，表示变量A1与自身的相关性是最大值（完全相关）。
- 横向和纵向的数值（例如0.90，0.96等）表明了变量之间的具体相关系数，数字越接近1，表示变量之间的正相关性越强。
- 可能在深红色区域中，一些变量之间显示出较强的相关性，提示可能存在某种联系或线性关系。

2、知识点感悟

1. 偏相关分析与相关性分析的区别

偏相关分析与相关性分析有着本质的区别。相关性分析（如Pearson相关）仅评估两个变量之间的线性关系，而偏相关分析则试图在控制其他变量的影响下理解两个变量之间的关系。在一些情况下，变量之间可能看似存在相关性，但实际上这种关系是由其他隐藏变量造成的。通过偏相关分析，可以更为准确地揭示变量间的因果关系，有助于做出更为科学的决策。

2. 数据处理的重要性

数据清洗是数据分析过程中的关键一步。在实际操作中，数据集常常会包含缺失值、异常值等，这些数据问题会直接影响到分析结果的准确性。通过 `dropna()` 等函数进行清洗，可以提高结果的可靠性。此外，对数据进行合适的重命名和格式化也有助于后续的分析工作，使代码更具可读性。

3. 可视化的重要性

绘制热力图是一种直观的方式，可以快速识别变量间的相关性和潜在的关系。这对于探索性数据分析尤为重要，能够帮助分析师发现数据中的模式、趋势以及异常。同时，合适的可视化也能够增强分析结果的表达力，使得非专业人士也能理解复杂的数据关系。

4. 计算工具的选择

在处理复杂数据分析时，选择合适的工具至关重要。在本例中，使用 `pandas` 进行数据处理，`seaborn` 和 `matplotlib` 进行可视化，以及 `pingouin` 进行统计分析，这些库都是高效且易于使用的，提高了开发效率，允许分析师将更多的精力投入到数据洞察中。

3、学习知识点的方法

1. 理论与实践相结合

理解统计分析的理论基础是学习的第一步。可以通过阅读相关书籍、在线课程等资源来掌握偏相关分析和Pearson相关分析的基本概念和理论。理论的学习应伴随实际的数据分析案例，可以选择公开数据集，如Kaggle、UCI Machine Learning Repository等，通过编程实现数据的清洗、统计分析和可视化。

2. 数据可视化工具的学习

除了案例分析，熟悉数据可视化工具的使用也是非常重要的。使用Python中的 `seaborn` 和 `matplotlib` 等库来进行数据可视化，通过实践来掌握如何绘制多种类型的图表，如热力图、散点图、柱状图等。

3. 参与社区与讨论

在学习的过程中，参与相关社区，如Stack Overflow、数据分析论坛等，可以帮助自己解决遇到的问题，同时与其他学习者和专业人士交流经验。这不仅能够扩展自己的知识面，还能够认识到数据分析的广泛应用。

4. 持续学习与更新

数据分析领域发展迅速，新工具、新方法层出不穷。因此，需要保持学习的热情，关注前沿技术和研究成果。通过观看在线讲座、参加工作坊等方式，能够及时获取最新的信息和技能。

4、经典案例

案例 1: 不同消费类型对销售额的影响分析

核心案例题目

分析某电商平台不同消费类型（如：电子产品、衣物、食品等）对总销售额的影响。

代码示例

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from pingouin import partial_corr

# 读取数据集
data = pd.read_csv('sales_data.csv')

# 处理缺失值
data_cleaned = data.dropna()

# 计算偏相关系数
variables = ['电子产品销售额', '衣物销售额', '食品销售额', '总销售额']
partial_corr_matrix = pd.DataFrame(index=variables, columns=variables)

for var1 in variables:
    for var2 in variables:
        if var1 != var2:
            result = partial_corr(data=data_cleaned, x=var1, y=var2, covar=variables)
            partial_corr_matrix.loc[var1, var2] = result['pearson'][0]
        else:
            partial_corr_matrix.loc[var1, var2] = 1.0

# 绘制热力图
plt.figure(figsize=(10, 8))
sns.heatmap(partial_corr_matrix.astype(float), annot=True, cmap='coolwarm')
plt.title('偏相关系数热力图')
plt.show()
```

代码说明：

1. 导入必要的库

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from pingouin import partial_corr
```

- `pandas`: 用于数据操作和分析, 尤其适合处理表格数据。
- `seaborn`: 基于Matplotlib而构建的可视化库, 提供更美观的统计图。
- `matplotlib.pyplot`: 提供绘图功能的库, 支持多种图表类型。
- `pingouin`: 一个用于统计分析的python库, 这里主要用于计算偏相关系数。

2. 读取数据集

```
data = pd.read_csv('sales_data.csv')
```

- 使用 `pd.read_csv` 读取名为 `sales_data.csv` 的CSV文件, 假定数据集中包含了销售额的相关信息。

3. 处理缺失值

```
data_cleaned = data.dropna()
```

- `dropna()` 函数用于删除数据帧中含有缺失值的行, 确保后续分析不受到缺失数据的影响。

4. 定义变量和准备偏相关系数矩阵

```
variables = ['电子产品销售额', '衣物销售额', '食品销售额', '总销售额']  
partial_corr_matrix = pd.DataFrame(index=variables, columns=variables)
```

- `variables` 是一个列表, 包含了要分析的变量名称。
- `partial_corr_matrix` 创建了一个空的Pandas数据框, 用于存放偏相关系数矩阵, 行索引和列索引都对应 `variables` 中的变量。

5. 计算偏相关系数

```
for var1 in variables:  
    for var2 in variables:  
        if var1 != var2:  
            result = partial_corr(data=data_cleaned, x=var1, y=var2, covar=variables)  
            partial_corr_matrix.loc[var1, var2] = result['pearson'][0]  
        else:  
            partial_corr_matrix.loc[var1, var2] = 1.0
```

- 外层循环遍历 `variables` 中的每个变量 `var1`。
- 内层循环遍历 `variables` 中的每个变量 `var2`。
- 如果 `var1` 和 `var2` 不同, 则调用 `partial_corr` 计算它们的偏相关系数, 并将结果存储在 `partial_corr_matrix` 中的相应位置。
- `covar=variables` 表示控制所有其他变量在进行偏相关分析时的影响。
- 如果 `var1` 和 `var2` 相同, 则相关系数固定为1.0, 因为任何变量与自身的相关性是完全的。

6. 绘制热力图

```
plt.figure(figsize=(10, 8))
sns.heatmap(partial_corr_matrix.astype(float), annot=True, cmap='coolwarm')
plt.title('偏相关系数热力图')
plt.show()
```

- `plt.figure(figsize=(10, 8))` 设置热力图的尺寸。
- `sns.heatmap(...)` 使用Seaborn库绘制热力图，`annot=True` 表示在单元格内显示数字，`cmap='coolwarm'` 是热图使用的颜色映射。
- 设置热力图标题为“偏相关系数热力图”。
- `plt.show()` 显示绘制好的热力图。

样例数据集：

示例数据集 `student_data.csv` 的表格，您可以根据这个表格创建相应的CSV文件。这个数据集包含四个变量：学习时间、上课参与度、课外活动和考试成绩。

学生ID	学习时间 (小时)	上课参与度 (%)	课外活动 (小时)	考试成绩 (%)
1	5	80	2	88
2	3	70	1	76
3	4	90	3	92
4	6	85	2	90
5	2	65	0	68
6	4	75	1	82
7	5	88	3	94
8	3	60	4	72
9	4	95	1	96
10	6	80	5	89

数据集说明：

- **学生ID**：唯一的学生识别码。
- **学习时间 (小时)**：每周用于学习的小时数。
- **上课参与度 (%)**：表示学生上课的出勤率，取值在0到100之间。
- **课外活动 (小时)**：参加课外活动的小时数。
- **考试成绩 (%)**：期末考试的分

结果解释

在热力图中可以清晰地看到各个消费类型之间的偏相关关系。例如，可能会发现“衣物销售额”与“总销售额”之间的偏相关系数比较高，说明在控制其他消费类型的影响后，这两个变量之间的相关性仍然很强，这可能意味着在促进衣物销售的同时，总销售额也有显著提升。

案例 2: 学生成绩与学习变量的关系分析

核心案例题目

分析学生的学习时间、上课参与度、课外活动对考试成绩的影响。

代码示例

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from pingouin import partial_corr

# 读取数据
data = pd.read_csv('student_data.csv')

# 清理缺失数据
data = data.dropna()

# 偏相关计算
variables = ['学习时间', '上课参与度', '课外活动', '考试成绩']
partial_corr_matrix = pd.DataFrame(index=variables, columns=variables)

for var1 in variables:
    for var2 in variables:
        if var1 != var2:
            result = partial_corr(data=data, x=var1, y=var2, covar=variables)
            partial_corr_matrix.loc[var1, var2] = result['pearson'][0]
        else:
            partial_corr_matrix.loc[var1, var2] = 1.0

# 绘制热力图
plt.figure(figsize=(10, 8))
sns.heatmap(partial_corr_matrix.astype(float), annot=True, cmap='coolwarm')
plt.title('学生成绩偏相关热力图')
plt.show()
```

结果解释

通过热力图，可以看出“学习时间”与“考试成绩”之间的偏相关系数非常高，可能说明学生投入更多学习时间会直接影响考试成绩的提升。而“课外活动”如果偏相关系数较低，则可能表明该因素对成绩的影响并不是直接的或显著的。

好的！让我们对这段代码进行详细说明，并生成一个示例数据集表格。

1. 导入必要的库

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from pingouin import partial_corr
```

- `pandas`: 用于数据处理和分析。
- `seaborn`: 用于绘制统计图表，提供更美观的可视化。
- `matplotlib.pyplot`: 提供绘图功能。
- `pingouin`: 用于进行统计分析，这里我们使用它来计算偏相关系数。

2. 读取数据

```
data = pd.read_csv('student_data.csv')
```

- 使用 `pd.read_csv` 函数从CSV文件中读取学生数据，数据文件名为 `student_data.csv`。

3. 清理缺失数据

```
data = data.dropna()
```

- `dropna()` 函数用于删掉含有缺失值的行，确保后面的分析可以依赖完整的数据集。

4. 定义变量和创建偏相关系数矩阵

```
variables = ['学习时间', '上课参与度', '课外活动', '考试成绩']  
partial_corr_matrix = pd.DataFrame(index=variables, columns=variables)
```

- `variables` 是一个列表，包含我们要分析的变量，即“学习时间”、“上课参与度”、“课外活动”和“考试成绩”。
- `partial_corr_matrix` 创建了一个空的 Pandas 数据框，其行和列均为 `variables`，用来存放偏相关系数。

5. 计算偏相关系数

```
for var1 in variables:  
    for var2 in variables:  
        if var1 != var2:  
            result = partial_corr(data=data, x=var1, y=var2, covar=variables)  
            partial_corr_matrix.loc[var1, var2] = result['pearson'][0]  
        else:  
            partial_corr_matrix.loc[var1, var2] = 1.0
```

- 这部分代码使用两个嵌套的循环遍历所有变量。
- 如果两个变量不同，调用 `partial_corr` 函数计算它们的偏相关系数，将结果存储到矩阵中；因为偏相关系数是控制其他变量的影响之后的相关性。
- 当两个变量相同时，相关系数被设置为 `1.0`，因为任何变量与自身的相关性是完全的。

6. 绘制热力图

```
plt.figure(figsize=(10, 8))  
sns.heatmap(partial_corr_matrix.astype(float), annot=True, cmap='coolwarm')  
plt.title('学生成绩偏相关热力图')  
plt.show()
```

- `plt.figure(figsize=(10, 8))` 设置热力图的尺寸。
- `sns.heatmap(...)` 使用Seaborn绘制热力图，`annot=True` 表示在图上显示数值，`cmap='coolwarm'` 指定了颜色映射。
- `plt.title` 设置热力图的标题为“学生成绩偏相关热力图”。
- `plt.show()` 用于显示绘制出的热力图。

示例数据集表格

以下是一个示例数据集 `student_data.csv` 的表格，包含四个变量——学习时间、上课参与度、课外活动和考试成绩。

学生ID	学习时间 (小时)	上课参与度 (%)	课外活动 (小时)	考试成绩 (%)
1	5	80	2	88
2	3	70	1	76
3	4	90	3	92
4	6	85	2	90
5	2	65	0	68
6	4	75	1	82
7	5	88	3	94
8	3	60	4	72
9	4	95	2	96
10	6	80	5	89

3.5 距离相关分析

本节案例完整的代码：

```
import pandas as pd
import matplotlib
matplotlib.use('TkAgg')
from scipy.spatial.distance import euclidean
from scipy.spatial.distance import cosine
# 1.导入数据
data = pd.read_csv('../data/ratings.csv')
# 2.根据userId、movieId、rating三列得到评分行
pivot_df = data.pivot(index='userId', columns='movieId', values='rating')
# 3.填充NaN值为0
pivot_df_subset = pivot_df.fillna(0)
print(pivot_df_subset)
# 4.创建新的DataFrame用于存储矩阵
matrix_df_a = pd.DataFrame(index=pivot_df_subset.index, columns=pivot_df_subset.index)
matrix_df_b = pd.DataFrame(index=pivot_df_subset.index, columns=pivot_df_subset.index)
# 5.填充欧式距离矩阵
for i in pivot_df_subset.index:
    for j in pivot_df_subset.index:
        # 计算用户i和用户j之间的欧式距离
        distance = euclidean(pivot_df_subset.loc[i], pivot_df_subset.loc[j])
        matrix_df_a.loc[i, j] = distance
print(matrix_df_a)
# 6.填充余弦相似度矩阵
for i in pivot_df_subset.index:
```

```

for j in pivot_df_subset.index:
    # 计算用户i和用户j之间的余弦相似度
    similarity = 1 - cosine(pivot_df_subset.loc[i], pivot_df_subset.loc[j])
    matrix_df_b.loc[i, j] = similarity
print(matrix_df_b)
# 7.保存结果
matrix_df_a.to_csv('距离相关分析-欧氏距离.csv')
matrix_df_b.to_csv('距离相关分析-余弦相似度.csv')

```

ratings.csv:

A	B	C	D	
userId	movieId	rating	timestamp	
1	1	4	964982703	
1	3	4	964981247	
1	6	4	964982224	
1	47	5	964983815	
1	50	5	964982931	
1	70	3	964982400	
1	101	5	964980868	
1	110	4	964982176	
1	151	5	964984041	
1	157	5	964984100	
1	163	5	964983650	
1	216	5	964981208	
1	223	3	964980985	
1	231	5	964981179	
1	235	4	964980908	
1	260	5	964981680	
1	296	3	964982967	
1	316	3	964982310	
1	333	5	964981179	
1	349	4	964982563	
1	356	4	964980962	
1	362	5	964982588	
1	367	4	964981710	
1	423	3	964982363	
1	441	4	964980868	
1	457	5	964981909	
1	480	4	964982346	
1	500	3	964981208	

610	152081	4	1.494E+09		
610	152372	3.5	1.494E+09		
610	155064	3.5	1.494E+09		
610	156371	5	1.48E+09		
610	156726	4.5	1.494E+09		
610	157296	4	1.494E+09		
610	158238	5	1.48E+09		
610	158721	3.5	1.48E+09		
610	158872	3.5	1.494E+09		
610	158956	3	1.494E+09		
610	159093	3	1.494E+09		
610	160080	3	1.494E+09		
610	160341	2.5	1.48E+09		
610	160527	4.5	1.48E+09		
610	160571	3	1.494E+09		
610	160836	3	1.494E+09		
610	161582	4	1.494E+09		
610	161634	4	1.494E+09		
610	162350	3.5	1.494E+09		
610	163937	3.5	1.494E+09		
610	163981	3.5	1.494E+09		
610	164179	5	1.494E+09		
610	166528	4	1.494E+09		
610	166534	4	1.494E+09		
610	168248	5	1.494E+09		
610	168250	5	1.494E+09		
610	168252	5	1.494E+09		
610	170875	3	1.494E+09		

为了确保提供详细和充实的内容报告，我需要逐步分析和提取图片中的信息，并结合代码的逻辑和实现。以下是按照要求生成的内容报告的框架：

1. 本节知识案例分析

这段代码是在使用Python进行数据分析，主要使用了 `pandas` 库来处理数据，`scipy` 库来计算不同用户之间的相似度和距离矩阵。整个过程分为几个步骤，具体如下：

1.1 导入数据

```
data = pd.read_csv('../data/ratings.csv')
```

首先，代码通过 `pandas` 导入了一个CSV文件，其中包含用户对电影的评分信息，包括用户ID、电影ID和评分值。

1.2 创建评分矩阵

```
pivot_df = data.pivot(index='userId', columns='movieId', values='rating')
```

接着，利用 `pivot` 函数将数据转化为一个透视表，其中行索引为用户ID，列索引为电影ID，值为评分。这样，评分矩阵就形成了，方便后续计算用户之间的相似度。

1.3 填充缺失值

```
pivot_df_subset = pivot_df.fillna(0)
```

缺失的评分（NaN值）被填充为0，这意味着对于用户未评分的电影，我们默认其评分为0，以避免在计算距离时出现问题。

1.4 创建距离矩阵

```
matrix_df_a = pd.DataFrame(index=pivot_df_subset.index, columns=pivot_df_subset.index)
matrix_df_b = pd.DataFrame(index=pivot_df_subset.index, columns=pivot_df_subset.index)
```

代码创建两个新的 `DataFrame` 来存储用户之间的距离和相似度。其中，`matrix_df_a` 用于存储欧氏距离，`matrix_df_b` 用于存储余弦相似度。

1.5 计算欧氏距离

```
for i in pivot_df_subset.index:
    for j in pivot_df_subset.index:
        distance = euclidean(pivot_df_subset.loc[i], pivot_df_subset.loc[j])
        matrix_df_a.loc[i, j] = distance
```

对于每一对用户，使用 `euclidean` 函数计算用户之间的欧氏距离，并将结果存储在矩阵中。欧氏距离反映了用户评分向量的整体差异。

1.6 计算余弦相似度

```
for i in pivot_df_subset.index:
    for j in pivot_df_subset.index:
        similarity = 1 - cosine(pivot_df_subset.loc[i], pivot_df_subset.loc[j])
        matrix_df_b.loc[i, j] = similarity
```

同样的，对于每一对用户，使用 `cosine` 函数计算其余弦相似度。余弦相似度反映了用户评分向量的方向相似性，值越接近1，表示用户之间的评分越相似。

1.7 保存结果

```
matrix_df_a.to_csv('距离相关分析-欧氏距离.csv')
matrix_df_b.to_csv('距离相关分析-余弦相似度.csv')
```

最后，分别将计算结果保存为CSV文件，方便后续分析和可视化。

[illegible]

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
	userId																				
1	1	0	0.02728258	0.05972020	0.19439470	0.12907965	0.12815155	0.15874444	0.13696808	0.06426336	0.01687525	0.13249905	0.01645755	0.09297135	0.11323830	0.16068875	0.16958584	0.26435794	0.21686474	0.32537610	0.1605
2	0	0.02728258	1	0	0.00372586	0.01661144	0.02533434	0.02758515	0.02725736	0	0.06744530	0.04441922	0	0.04391833	0.01690122	0.19777555	0.03972755	0.16625307	0.01257055	0.0187	
3	0	0.05972020	0	0	0.00223130	0.00501970	0.03993635	0	0.00494115	0	0	0	0	0	0.00360387	0.01725142	0.03299215	0.00893110	0.02824005	0.01914155	0.0041
4	0.19439470	0.00372586	0.02253138	1	0.12656589	0.08849102	0.11511990	0.06269690	0.11360665	0.03163335	0.05476677	0.04994477	0.07694858	0.04988866	0.11505077	0.16476087	0.14568303	0.12312637	0.20665300	0.1164	
5	0.12907965	0.01661144	0.00501970	0.08849102	1	0.30034895	0.10834170	0.42907545	0	0.03610695	0.18380522	0.05886021	0.07156565	0.22113136	0.01151705	0.08217085	0.06233005	0.12131335	0.09875840	0.0934	
6	0.12815155	0.02533434	0.03993635	0.08849102	0.30034895	1	0.07584315	0.37048815	0.01394055	0.02038465	0.22023485	0.09605538	0.02391775	0.34644245	0.06756605	0.04927485	0.06061915	0.11582802	0.22186870	0.0855	
7	0.15874444	0.02758515	0	0	0.11511990	0.10834170	0.07584315	1	0.11488495	0.09462626	0.14967805	0.04233465	0.02309075	0.04736165	0.23047505	0.15879008	0.06261915	0.12728210	0.12696805	0.1777	
8	0.13696808	0.02757300	0.00494115	0.06269690	0.42907545	0.37048815	0.11488495	1	0	0.02391420	0.23533370	0.04619465	0.05441715	0.43006411	0.12819595	0.14719805	0.09996110	0.15217245	0.16918240	0.0513	
9	0.06426336	0	0	0.01136065	0	0.01394055	0.01494626	0	1	0.04059591	0	0.05021725	0	0.12707534	0.09228801	0.08991110	0.10667895	0.04886955	0.0937		
10	0.01687525	0.06744530	0	0.03116335	0.03610695	0.02038465	0.13029945	0.02391420	0.04059591	1	0.02727985	0.09115975	0.03185471	0.02188592	0.05233675	0.09818955	0.02357878	0.15139005	0.06799875	0.0805	
11	0.13249905	0.04441922	0	0.05476677	0.18380522	0.22023485	0.19467805	0.23533370	1	0.02727985	1	0.03189822	0.03057915	0.02788905	0.09052225	0.05384477	0.17108885	0.15034455	0.13660231		
12	0.01645755	0	0	0.04994477	0.05886020	0.09603585	0.04523465	0.04619465	0	0.09115975	0.03189822	1	0.03784627	0.01949335	0.04265377	0	0.02698305	0.01358875	0.07366125	0.0267	
13	0.09297135	0.03918335	0	0.07694858	0.07156565	0.02391775	0.02380450	0.04441715	0.05021725	0.03814645	0.03057915	0.03748627	1	0.06593272	0.16074277	0.08624485	0.07394545	0.09737327	0.06941185	0.0922	
14	0.11323830	0.01690122	0.00360387	0.04988866	0.22171313	0.34644245	0.04736165	0.05046115	0	0.02718522	0.08079995	0.01949335	0.06593272	1	0.07679522	0.03859845	0.13731545	0.11860575	0.09891745	0.0062	

学习数据分析与相似度计算相关知识点的方式有很多，以下是一些有效的方法：

3.1 阅读相关书籍

通过阅读数据分析、机器学习和推荐系统的书籍，可以系统地学习相关理论和算法。例如，《集体智慧编程》介绍了如何用Python编写推荐系统。

3.2 在线课程

借助平台如Coursera、Udacity等，可以找到数据分析和机器学习相关的在线课程，学习数据处理、分析和建模的实用技巧。

3.3 实践项目

自行创建一个简单的推荐系统项目，应用所学知识，通过真实数据的练习加深理解。例如，利用MovieLens数据集进行用户评分的分析和可视化。

3.4 社区交流

加入相关技术社区，如Stack Overflow、Kaggle等，向他人请教问题，分享自己的理解与实践经验，从而促进学习。

3.5 编写总结报告

对学习内容进行总结，包括代码实现、应用场景、自己的理解和见解等，这不仅帮助记忆，也能加深对知识的理解。

4. 经典案例

结合以上所提到的知识点，我们可以延伸出以下经典案例：

4.1 案例1: 用户基于内容的推荐系统

案例标题：基于相似用户的推荐系统

```
import pandas as pd
from scipy.spatial.distance import euclidean, cosine

# 加载数据
data = pd.read_csv('ratings.csv')

# 创建评分矩阵
pivot_df = data.pivot(index='userId', columns='movieId', values='rating').fillna(0)

# 计算用户相似度
user_similarity = pd.DataFrame(index=pivot_df.index, columns=pivot_df.index)

for u1 in pivot_df.index:
    for u2 in pivot_df.index:
        user_similarity.loc[u1, u2] = 1 - cosine(pivot_df.loc[u1], pivot_df.loc[u2])

# 输出相似度矩阵
print(user_similarity)
```

在这个案例中，首先通过 pandas 加载评分数据并创建评分矩阵。然后，利用余弦相似度计算出用户之间的相似度，并将结果输出。这个案例强调了用户相似性在推荐中的应用。

数据集表格示例

userId	movieId	rating
1	101	5
1	102	3
1	103	4
2	101	2
2	102	5
2	104	4
3	101	4
3	103	5
4	102	3
4	104	2

数据集说明

- **userId**: 用户唯一标识符，表示该评分记录是哪个用户的评分。
- **movieId**: 电影唯一标识符，表示该评分记录对应的电影。
- **rating**: 用户对电影的评分（通常为1到5的整数），反映用户对电影的喜好程度。

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	userid	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
2	1	0.02728286	0.05972026	0.19439476	0.12907985	0.12815157	0.15874444	0.13696806	0.06426336	0.01687525	0.13249906	0.01645755	0.09297133	0.11323836	0.16068876	0.16985846	0.26435794	0.21486847	0.32537616	0.1605	
3	2	0.02728286	1	0	0.00372586	0.01661444	0.02533345	0.02758515	0.02725736	0	0.06744536	0.04441925	0	0.04391837	0.01690127	0.11977755	0.09372751	0.10375496	0.16625307	0.01257057	0.0141
4	3	0.05972026	0	1	0.00225136	0.00501971	0.00393631	0	0.00494111	0	0	0	0	0.00306382	0.01725147	0.03229916	0.00981311	0.02824056	0.01914156	0.014156	0.0087
5	4	0.19439476	0.00372586	0.00225136	1	0.12658599	0.08849105	0.11511991	0.06296891	0.01136067	0.03116335	0.05476677	0.04994471	0.07694888	0.04898867	0.07155077	0.16476087	0.14505835	0.12321677	0.20605307	0.1137
6	5	0.12907985	0.01661444	0.00501971	0.12658599	1	0.30034897	0.10834176	0.42907545	0	0.03061096	0.18380522	0.05886021	0.01715658	0.22171136	0.11015175	0.08217085	0.16263306	0.12131336	0.09875841	0.0964
7	6	0.12815157	0.02533345	0.00393631	0.08849105	0.30034897	1	0.07584316	0.37048816	0.01390405	0.02038465	0.22023485	0.09605588	0.02391776	0.34644245	0.06757605	0.04927485	0.09068197	0.11558206	0.22184877	0.0855
8	7	0.15874444	0.02758515	0	0.11511991	0.10834176	0.07584316	1	0.11488495	0.09946266	0.13209945	0.19467807	0.04523466	0.03208046	0.07473615	0.23047506	0.15879086	0.24650156	0.27282816	0.12696807	0.1727
9	8	0.13696806	0.02757361	0.00494111	0.06296891	0.42907545	0.37048816	0.11488495	1	0	0.02391427	0.23533371	0.04619465	0.05441714	0.43006411	0.10281956	0.11471906	0.19999646	0.15217247	0.16918247	0.0513
10	9	0.06426336	0	0	0.11136067	0	0.01390405	0.09946266	0	1	0.04095914	0	0	0.05021726	0	0.12707534	0.07922801	0.08981111	0.11066986	0.04886957	0.0937
11	10	0.01687525	0.06744536	0	0.03116335	0.03061096	0.02038465	0.13209945	0.02391427	0.04095914	1	0.02727988	0.09115976	0.03814641	0.02718525	0.10825367	0.09818955	0.07235786	0.11513907	0.02679987	0.0805
12	11	0.13249906	0.04441925	0	0.05476677	0.18380522	0.22023485	0.19467807	0.02353371	0	0.02727988	1	0.03189822	0.03057911	0.20789905	0.09025226	0.05384477	0.17108887	0.15030456	0.13960231	
13	12	0.01645755	0	0	0.04994471	0.05886021	0.09605588	0.04523466	0.04619465	0	0.09115976	0.03189822	1	0.03784627	0.01941936	0.04265376	0	0.02698304	0.01358877	0.07366127	0.0267
14	13	0.09297133	0.04391837	0	0.07694888	0.07156587	0.02391776	0.03208046	0.05441714	0.05021726	0.03814641	0.03057911	0.03784627	1	0.06593271	0.10674271	0.08624488	0.07398454	0.07973277	0.10644187	0.0928
15	14	0.11323836	0.1690127	0.00306382	0.04898867	0.22171136	0.34644245	0.07473615	0.43006411	0	0.02718525	0.20789905	0.01941936	0.06593271	1	0.07679522	0.08358945	0.16731546	0.11860576	0.08991747	0.0062
16	15	0.16068876	0.11977755	0.01725147	0.07155077	0.11015175	0.06757605	0.23047506	0.10281956	0.12707534	0.08253671	0.09025226	0.04265376	0.10674271	0.07679522	1	0.23167336	0.30507023	0.30697137	0.15777324	0.1306
17	16	0.16985846	0.09372751	0.03229916	0.16476087	0.08217085	0.04927485	0.15879086	0.11471906	0.07922801	0.09818955	0.05384477	0	0.08624488	0.08358945	0.23167336	1	0.45609561	0.30365097	0.15048961	0.0426
18	17	0.26435794	0.10375496	0.00981311	0.14505835	0.16263306	0.09068197	0.24650156	0.19999646	0.08981111	0.07235786	0.17108887	0.02698304	0.07398454	0.16731546	0.30507023	0.45609561	1	0.32121246	0.17819387	0.0534
19	18	0.21486847	0.16625307	0.02824056	0.12321677	0.12131336	0.11558206	0.27282816	0.15217247	0.11066986	0.11513907	0.15030456	0.01358877	0.07973277	0.11860576	0.30697137	0.30365097	0.32121246	1	0.18932185	0.0943
20	19	0.32537616	0.01257057	0.01914156	0.20605307	0.09875841	0.22184877	0.12696807	0.16918247	0.04886957	0.02679987	0.13960231	0.07366127	0.10644187	0.08991747	0.15777324	0.15048961	0.17819387	0.18932185	1	0.2023
21	20	0.16096937	0.01413695	0.00875594	0.11375485	0.09647435	0.08554375	0.17279977	0.05134123	0.09375406	0.08051641	0	0.02673355	0.09286114	0.00624211	0.13089366	0.04261334	0.05341227	0.09433704	0.20230334	
22	21	0.15316181	0.09087976	0.00401674	0.05301394	0.05826386	0.06252185	0.25762437	0.07217021	0.10633725	0.21029285	0.09972277	0.04275491	0.06864296	0.04112015	0.28800027	0.11463271	0.15219365	0.32910955	0.15565991	0.1139
23	22	0.05069055	0.14463514	0.00306955	0.04340306	0.03307395	0.04178535	0.13701597	0.07234736	0.05485735	0.13259027	0.04421208	0	0.07212747	0.04361376	0.18145967	0.16277076	0.16533945	0.21017345	0.06527047	0.1106
24	23	0.10666935	0.01359733	0.00308111	0.09397141	0.06688884	0.02834755	0.08105647	0.05277047	0.02729995	0.02117505	0.01577904	0	0.02047656	0.02451577	0.08582584	0.32124602	0.23542247	0.20028107	0.07769855	0.0373
25	24	0.15519346	0.12990177	0.00259247	0.07169033	0.09621515	0.08234544	0.11677097	0.13086185	0.03761005	0.15605405	0.17093625	0.02695847	0.09168321	0.09050915	0.22712596	0.17249033	0.18812217	0.25475733	0.10496796	0.0431
26	25	0.09926134	0.22600833	0.00510005	0.03286196	0.04070516	0.01723677	0.09212314	0.06571171	0.12867525	0.14809905	0.03264798	0	0.14525925	0.04472055	0.23926505	0.24281797	0.25276807	0.19003724	0.04878347	0.0502
27	26	0.10579148	0	0	0.06183224	0.29428176	0.24324216	0.13154095	0.44886747	0	0.02331095	0.29557078	0	0.06307296	0.30205911	0.07069206	0.08618331	0.14862567	0.10620356	0.10922094	0.0135
28	27	0.23893291	0	0.00233165	0.12895386	0.10049125	0.11090745	0.15793735	0.11357311	0.06471017	0.02228706	0.12089985	0.06557941	0.05644788	0.05773885	0.12598206	0.12192494	0.14311234	0.13662488	0.24235166	0.2107
29	28	0.20286555	0.05814521	0.02977117	0.13559325	0.12189448	0.10242926	0.23970676	0.10703667	0.07214277	0.10102325	0.13157574	0.01492522	0.08910566	0.10870384	0.26043234	0.22350691	0.27355257	0.33986366	0.16308089	0.0873

该数据集用于构建评分矩阵，从而计算用户之间的相似度，进而在推荐系统中实现用户推荐功能。

4.2 案例 2: 商品推荐

案例标题：基于相似商品的推荐系统


```
import pandas as pd
from scipy.spatial.distance import cosine

# 加载商品评分数据
data = pd.read_csv('item_ratings.csv')

# 创建商品评分矩阵
pivot_df = data.pivot(index='userId', columns='itemId', values='rating').fillna(0)

# 计算商品之间的相似度
item_similarity = pd.DataFrame(index=pivot_df.columns, columns=pivot_df.columns)

for item1 in pivot_df.columns:
    for item2 in pivot_df.columns:
        item_similarity.loc[item1, item2] = 1 - cosine(pivot_df[item1], pivot_df[item2])

# 输出商品相似度矩阵
print(item_similarity)
```

在这个商品推荐的案例中，代码的逻辑类似于用户推荐系统，但是主要计算的是商品之间的相似度。这表明如何通过相似商品来进行推荐，增加用户的购买可能性。

数据集表格示例

userId	itemId	rating
1	A	5
1	B	3
2	A	4
2	C	2
2	B	5
3	A	2
3	C	4
4	B	5
4	C	1

数据集说明

- **userId**: 用户唯一标识符，表示该评分记录是哪个用户的评分。
- **itemId**: 商品唯一标识符，表示该评分记录对应的商品。
- **rating**: 用户对商品的评分（通常为1到5的整数），反映用户对商品的喜好程度。

该数据集用于构建商品评分矩阵，从而计算商品之间的相似度，帮助实现根据用户喜好进行商品推荐的功能。

输出结果解释

输出的结果包含距离矩阵和相似度矩阵，这两个矩阵可以帮助我们理解用户之间或商品之间的相似性。例如，用户1与用户2之间的余弦相似度 high 说明他们的评分特点相似，可能适合向用户1推荐用户2喜欢的商品。

5. 学习效果检验

1. 题目的难度和目的

这组题目主要集中在统计学与数据分析的领域，尤其是相关性分析与回归分析。题目的整体难度较高，适合掌握了一定统计基础知识的学生或学习者。在题目中，涉及到使用不同的相关系数（如 Pearson、Spearman 和 Kendall 等）来评估变量之间的关系，这一内容通常需要具备一定的数学背景，能够理解和运用相应的公式与图表。

2. 题目涉及的知识点

1. 相关系数的定义与计算：

- Pearson 相关系数：适用于衡量两个连续变量之间的线性关系，需要满足正态分布的假设。
- Spearman 相关系数：用于评估两个变量的单调关系，尤其在数据不满足正态分布时，也可以较为有效地进行分析。
- Kendall 相关系数：这是一种比较两个变量秩次（rank）一致性的测量，适用于小样本或者包含极端值的情况。

2. 假设检验：

- 在统计学中，假设检验是非常重要的部分，考生需要熟悉如何进行构造、检验和判断假设。配合相关系数的计算，考生需要了解 p 值的含义，如何通过 p 值判断相关性是否显著。

3. 数据的描述性统计：

- 在数据分析中，描述性统计是基础，涉及均值、中位数、标准差等概念。考生需理解这些统计量的含义及其在相关分析中的重要性。

4. 使用编程语言实现统计分析：

- 题目还涉及如何利用 Python 等编程语言进行统计分析，要求考生能够将理论知识转化为代码，熟练运用相关库（如 pandas、numpy、scipy 等）完成计算与图表的生成。

3. 对题目的感想

从整体来看，这组题目挑战性较大，特别是在涉及到理论知识的应用时。不过，这也反映了统计学在实际工作中的复杂性与重要性。统计学不仅是一门理论学科，更是实践中的工具，能够帮助我们从大量数据中提取有用的信息。通过学习与解题，我深刻体会到以下几点：

1. 实践与理论的结合：

- 在学习统计学时，单纯的理论学习往往无法全面掌握知识。实际的案例分析和编程实践能够有效提高理解和使用能力，因此在备考或学习时，应更加注重实践。

2. 统计工具的重要性：

- 学会运用 Python 等编程工具不仅能提高效率，还能提升数据分析的准确性。在未来的学习或工作中，熟练掌握这些工具将是提升自我竞争力的重要因素。

3. 细致耐心的重要性：

- 在相关分析中，数据的细节往往能够影响结果及结论。因此，在进行数据分析时，保持细致和耐心，认真对待每一个步骤是至关重要的。

4. 举例说明

关于统计学中Spearman 相关系数的题目来进行分析并给出解决方法。假设题目内容如下：

题目描述

给定一组数据，包含两个变量：X 和 Y。请计算 X 和 Y 之间的 Spearman 相关系数，并判断它们之间的相关性强弱。

数据如下：

X	Y
1	2
2	1
3	4
4	3
5	5

解决方法

1. 理论知识回顾

Spearman 相关系数用于测量两个变量的单调关系，尤其适用于非正态分布数据。计算步骤通常包括：

- 对每个变量的数值进行排序，并为其分配秩位。
- 计算每对秩位的差。
- 使用 Spearman 的公式 ($\rho = 1 - \frac{\sum d_i^2}{n(n^2 - 1)}$) 计算相关系数，其中 (d_i) 是每对秩位的差， n 是数据点的数量。

2. 计算过程

首先，将给定的 X 和 Y 排序并分配秩位：

X	Rank_X	Y	Rank_Y
1	1	2	2
2	2	1	1
3	3	4	4
4	4	3	3
5	5	5	5

然后，计算每对秩位的差 (d_i):

Rank_X	Rank_Y	(d_i)	(d_i^2)
1	2	-1	1
2	1	1	1
3	4	-1	1
4	3	1	1
5	5	0	0

接下来，计算 (\sum d_i^2 = 1 + 1 + 1 + 1 + 0 = 4)。n = 5。

使用 Spearman 的公式计算相关性：

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad \rho = 1 - \frac{n(n^2 - 1)6 \sum d_i^2}{n(n^2 - 1)}$$

3. 结果分析

Spearman 相关系数 (= 0.8) 表示变量 X 和 Y 之间具有很强的正相关性。

4. 编程实现 (Python 示例)

若你希望使用编程工具进行计算，可以使用 Python 的 `scipy` 库来简化这个过程：

```
import numpy as np
from scipy.stats import spearmanr

# 定义数据
X = np.array([1, 2, 3, 4, 5])
Y = np.array([2, 1, 4, 3, 5])

# 计算 Spearman 相关系数
correlation, p_value = spearmanr(X, Y)

print(f"Spearman 相关系数: {correlation:.2f}")
print(f"p值: {p_value:.4f}")
```

总结

通过以上分析，我们成功计算出了 X 和 Y 之间的 Spearman 相关系数，并得出了它们之间的强正相关性。这道题目考察了对非参数检验的理解与应用，如果对其他相关统计概念有疑问，请随时问我！

6. 选择题和客观题解答过程

选择题

1. 相关分析研究的是 () 之间的关系。

- A. 因果关系
- B. 相互影响关系
- C. 相关性关系
- D. 相似性

答案：C. 相关性关系

解释：

相关分析旨在探讨两个或多个变量之间的相关性，即它们在统计意义上的关联程度。相关关系描述了变量之间是否存在某种联系，以及这种联系的强度和方向，但并不涉及因果关系的确定。选项C直接指出了“相关性关系”，最符合相关分析的定义。

2. Pearson相关系数的取值范围是（ ）。

- A. $[-1, 1]$
- B. $[0, \infty)$
- C. $[-\infty, \infty)$
- D. $[0, \infty]$

答案：A. $[-1, 1]$

解释：

Pearson相关系数衡量两个连续变量之间的线性关系，其取值范围介于-1到1之间。值为1表示完全正相关，值为-1表示完全负相关，值为0表示无线性相关关系。选项A准确描述了这一范围。

3. 傅里叶分析的作用是（ ）。

- A. 研究变量之间的相关关系
- B. 提供基本的相关性计算条件
- C. 评估相关性的相似性
- D. 确定变量间的线性关系

答案：D. 确定变量间的线性关系

解释：

傅里叶分析主要用于将信号分解为不同频率的正弦和余弦成分，广泛应用于信号处理、图像分析等领域。虽然选项D“确定变量间的线性关系”不是傅里叶分析的直接定义，但在信号处理中的应用中，傅里叶分析帮助理解信号的频率成分，这可以间接用于分析变量之间的线性特性。因此，选项D最为接近。

4. 欧氏距离分析，欧氏距离和曼哈顿距离的区别是（ ）。

- A. 欧氏距离考虑个体的数量，曼哈顿距离考虑个体的密度
- B. 欧氏距离考虑每个维度的距离，曼哈顿距离考虑每个维度的数值
- C. 欧氏距离是直线距离，曼哈顿距离是沿坐标轴的距离
- D. 欧氏距离是距离的平方和，曼哈顿距离是绝对值和

答案：C. 欧氏距离是直线距离，曼哈顿距离是沿坐标轴的距离

解释：

欧氏距离（Euclidean Distance）是指在欧几里得空间中两点之间的直线距离，计算方法基于勾股定理。曼哈顿距离（Manhattan Distance）则是指在城市街区式布局中，两点之间沿坐标轴的距离总和，类似于走路时只能沿街道方向移动的距离。选项C准确描述了两者的主要区别。

5. 在显著性检验中，使用（ ）评价样本间的相关性程度。

- A. Pearson相关系数
- B. Spearman相关系数
- C. Kendall相关系数
- D. 距离度量方法

答案：A. Pearson相关系数

解释：

在显著性检验中，Pearson相关系数常用于评估两个连续变量之间的线性相关性程度。通过计算Pearson相关系数并进行显著性检验，可以判断观察到的相关性是否具有统计学意义。虽然Spearman和Kendall相关系数也可用于相关性分析，但Pearson相关系数是最常用于显著性检验的标准方法。故选A。

6. 以下（ ）相关系数适用于衡量分类变量之间的相关性。

- A. Pearson相关系数
- B. Kendall等级相关系数
- C. Spearman等级相关系数
- D. 距离相关系数

答案：D. 距离相关系数

解释：

传统的相关系数如Pearson、Spearman和Kendall主要用于连续或有序变量之间的相关性分析。对于分类变量之间的相关性，距离相关系数（Distance Correlation）更为适用，因为它能够捕捉到变量之间的复杂关系，包括非线性关系。选项D最符合衡量分类变量相关性的需求。

7. 数据集存在异常值时，使用（ ）相关系数更合适。

- A. Pearson相关系数
- B. Spearman相关系数
- C. Kendall相关系数
- D. 方差相关系数

答案：B. Spearman相关系数

解释：

当数据集中存在异常值时，Pearson相关系数可能会受到极端值的显著影响，导致结果不可靠。Spearman相关系数基于秩次（排名）计算，不受异常值的影响，因而在处理含有异常值的数据时更为稳健和合适。虽然Kendall相关系数也具有稳健性，但Spearman相关系数更常用，因此选B。

8. 在数据分析中，以下（ ）相关系数更适合用于计算样本之间的非线性特性。

- A. Pearson相关系数
- B. Spearman相关系数
- C. Kendall相关系数
- D. 非参数统计方法

答案：D. 非参数统计方法

解释：

非参数统计方法不依赖于数据的具体分布形式，适用于分析样本之间的复杂关系，包括非线性特性。相比之下，Pearson相关系数仅能捕捉线性关系，而Spearman和Kendall相关系数虽然可以处理某些非线性单调关系，但非参数方法提供了更广泛的适用性和灵活性。因此，选项D是最合适的。

9. Kendall相关系数相对于Spearman等相关系数的优点是（ ）。

- A. 运算是简单点
- B. 结果更稳健
- C. 计算效率低
- D. 相关结果更小

答案：B. 结果更稳健

解释：

Kendall相关系数（Kendall's Tau）比Spearman相关系数在处理数据中的小样本和存在并列排名的情况下更为稳健。它对数据中的异常值和分布不均情况具有更高的鲁棒性，能够提供更可靠的相关性评估。选项B准确反映了Kendall相关系数的主要优势。

判断题

1. 相关系数只能用于衡量线性关系，不能用于衡量非线性关系。（错误）

解释： 相关系数，尤其是皮尔逊相关系数，主要用于衡量线性关系。然而，并非只能用于线性关系。其他相关系数如斯皮尔曼等级相关系数和距离相关系数可以用于衡量非线性关系。因此，原题的表述过于绝对，实际情况更加复杂。

2. 偏相关系数的值越大，表示两个变量之间的关系越强。（正确）

解释： 偏相关系数衡量在控制其他变量影响后，两个变量之间的线性关系强度。其绝对值越大，表示在排除其他变量影响后，这两个变量之间的关系越强。因此，偏相关系数的值越大，关系确实越强。

3. 在相关关系中，控制变量对其他变量进行控制，保持其不变。（正确）

解释： 控制变量的目的是排除其对主要变量关系的影响，从而更准确地评估主要变量之间的相关性。在相关分析中，通过控制变量，可以保持其不变，专注于研究两个主要变量之间的关系。

4. 若两个变量间的相关系数大于0.6，则它们的相关性比相关系数大于0.8的重要程度。（错误）

解释： 相关系数的绝对值越大，表示两个变量之间的相关性越强。题目中“相关系数大于0.6”相比“相关系数大于0.8”具有较弱的相关性，因此说其“重要程度”更高是不正确的。

5. 距离相关分析中，欧氏距离的计算方法是将两点之间的平方和开根号。（正确）

解释： 欧氏距离是最常用的距离度量方法，其计算公式确实是两点之间各维度差值的平方和再开平方根。因此，题目的描述是正确的。

6. 相关系数的表征两个变量之间是否存在任何关系。（错误）

解释： 相关系数主要衡量两个变量之间线性关系的强度和方向。它并不能全面表征两变量之间存在的所有类型关系，尤其是非线性关系。因此，说相关系数表征“任何”关系是不准确的。

7. 相关关系可以用于判断因果关系。（错误）

解释： 相关关系只表明两个变量之间存在统计上的关联，但并不能确定因谁因果关系。因果关系的确定需要更深入的研究设计和分析，如实验设计或因果推断方法。

8. 如果两个变量之间的相关关系大于0.9，则可以确定它们之间存在直接的依赖关系。（错误）

解释： 虽然高度相关可能暗示强关联，但这并不意味着存在直接的因果依赖关系。高相关性可能由第三方变量影响或其他复杂因素引起，因此不能仅凭相关系数大于0.9就确定直接依赖关系。

9. 相关性分析可用于判断两个变量之间的相关性。（正确）

解释：相关性分析的主要目的就是评估和判断两个变量之间的相关性，包括相关性的方向和强度。因此，题目的描述是正确的。

10. Spearman等级相关系数可以用于判断两个变量之间的线性关系。（错误）

解释：斯皮尔曼等级相关系数用于评估两个变量之间的单调关系，不限于线性关系。它更适用于非参数数据或当数据不满足线性相关的假设时使用。因此，其主要用于判断单调而非严格的线性关系。

简答题

1. 简述相关系数的定义，并列举常用的相关系数及其适用场景。

回答：

相关系数的定义：

相关系数（Correlation Coefficient）是用于衡量两个变量之间线性关系强度和方向的统计指标。其取值范围通常在-1到1之间，正值表示正相关，负值表示负相关，绝对值越接近1表示相关性越强，接近0则表示相关性较弱或无相关性。

常用的相关系数及其适用场景：

1. 皮尔逊相关系数（Pearson Correlation Coefficient）：

- **定义：** 衡量两个连续变量之间线性关系的强度和方向。
- **适用场景：** 数据满足正态分布，适用于连续型数据。例如，研究身高与体重之间的关系。

2. 斯皮尔曼等级相关系数（Spearman's Rank Correlation Coefficient）：

- **定义：** 基于变量的排名来衡量两个变量之间的单调关系。
- **适用场景：** 数据不满足正态分布或存在异常值，适用于序数型或非线性关系的连续型数据。例如，分析学生成绩排名与学术表现的关系。

3. 肯德尔秩相关系数（Kendall's Tau）：

- **定义：** 基于变量排名的一致性来衡量两个变量之间的相关性。
- **适用场景：** 数据量较小或存在大量相同排名时使用，适用于序数型数据。例如，评估不同评审者对项目的排名一致性。

4. 点二列相关系数（Point-Biserial Correlation Coefficient）：

- **定义：** 衡量一个二元变量和一个连续变量之间的相关性。
- **适用场景：** 一个变量是二分类的，另一个是连续型的。例如，研究性别（男/女）与薪资水平之间的关系。

5. Phi系数（ Φ ）：

- **定义：** 衡量两个二元变量之间的相关性。
- **适用场景：** 适用于两个二分类变量。例如，分析吸烟与患某疾病的关联性。

这些相关系数根据数据类型和分布特征的不同，提供了多样化的方法来评估变量间的关系，帮助研究者选择最适合其数据特性的统计工具。

2. 相关系数可以用来判断两个变量之间的因果关系吗？为什么？

回答：

相关系数本身**不能**用来判断两个变量之间的因果关系。原因如下：

1. **相关不等于因果：** 高相关性仅表明两个变量在统计上存在关联，但并不意味着一个变量的变化会引起另一个变量的变化。可能存在第三个变量同时影响这两个变量，导致它们表现出相关性。

2. **缺乏时间序列信息：** 因果关系通常需要确定变量之间的时间顺序，即原因发生在结果之前。相关系数无法提供这种时间信息。
3. **潜在混杂因素：** 两个变量之间的相关性可能是由于其他未被考虑的变量（混杂因素）所导致。例如，冰淇淋销量的增加与溺水事件增多之间可能存在正相关，但实际的原因是夏季气温升高。
4. **非线性关系可能被误判：** 相关系数主要衡量线性关系，对于非线性因果关系，相关系数可能无法准确反映变量之间的联系。

结论：

要确定因果关系，通常需要通过实验设计（如随机对照试验）、时间序列分析或因果推断方法（如格兰杰因果检验、工具变量法）等更为复杂的统计方法，而不仅仅依赖于相关系数。

3. 相关系数的巨变是否受到样本大小的影响？为什么？

回答：

相关系数的波动确实受到样本大小的影响。样本大小对相关系数的估计稳定性和显著性有显著影响，具体原因如下：

1. 估计稳定性：

- **小样本：** 在样本量较小的情况下，相关系数的估计值可能具有较大的波动性和不确定性，容易受到异常值或单个观测点的影响，导致估计结果不稳定。
- **大样本：** 随着样本量的增加，相关系数的估计趋于稳定，减少了由于偶然性引起的偏差，提高了估计的准确性和可靠性。

2. 统计显著性：

- **小样本：** 检验相关系数是否显著时，小样本可能导致功效不足，即使存在实际的相关性，也可能因为样本量不足而无法检测到显著性。
- **大样本：** 大样本提高了检验的统计功效，使得即使是微小的相关性也可能达到统计显著性。

3. 置信区间：

- **小样本：** 相关系数的置信区间较宽，反映了较大的不确定性。
- **大样本：** 置信区间较窄，表示更高的估计精度。

总结：

样本大小对相关系数的估计和其解释具有重要影响。较大的样本能够提供更稳定、更可靠的相关系数估计，而较小的样本可能导致相关系数的不稳定和不准确。因此，在进行相关性分析时，应考虑样本量的充分性，以确保结果的有效性和可靠性。

4. 什么是偏相关系数？它与普通相关系数有何区别？

回答：

偏相关系数 (Partial Correlation Coefficient)：

偏相关系数用于衡量在控制了其他变量影响后，两个变量之间的线性关系强度和方向。具体来说，它度量的是两个变量在排除其他变量影响后的直接相关性。

与普通相关系数的区别：

1. 控制变量的存在：

- **普通相关系数：** 衡量两个变量之间的总体线性关系，不考虑其他潜在变量的影响。
- **偏相关系数：** 控制了一个或多个其他变量后，衡量两个变量之间的直接线性关系。

2. 反映关系的深度：

- **普通相关系数：** 可能受到第三方变量的影响，无法区分直接关系与间接关系。

- **偏相关系数**：通过控制其他变量，能够更准确地反映两个变量之间的直接关系。

3. 应用场景：

- **普通相关系数**：适用于简单的双变量分析。
- **偏相关系数**：适用于多变量分析，尤其是在需要理解变量间纯粹关系时，如多元回归分析中的独立变量间的相关性评估。

举例说明：

假设研究教育水平（X）与收入（Y）之间的关系，同时控制了工作经验（Z）的影响。

- **普通相关系数（ r_{XY} ）**：衡量教育水平与收入之间的总体相关性，未控制工作经验。
- **偏相关系数（ $r_{XY.Z}$ ）**：衡量教育水平与收入之间在控制了工作经验后的直接相关性。

总结：

偏相关系数通过控制其他变量的影响，提供了对两个变量之间纯粹相关性的更深入理解，相较于普通相关系数，它能揭示更为精确的变量关系模式。

5. 欧氏距离和曼哈顿距离有何区别？

回答：

欧氏距离（Euclidean Distance）和**曼哈顿距离（Manhattan Distance）**都是衡量两点之间距离的常用方法，但它们的计算方式和适用场景有所不同。

1. 定义与计算方式：

• 欧氏距离：

- **定义**：在欧几里得空间中，两点之间的“直线”距离。
- **计算公式（二维空间）**：

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

- **几何解释**：两点之间的最短路径，即直线距离。

• 曼哈顿距离：

- **定义**：在网格状的城市街道中，两点之间沿着轴对齐的路径的总距离。
- **计算公式（二维空间）**：

$$d = |x_2 - x_1| + |y_2 - y_1|$$

- **几何解释**：沿着轴方向（如城市街道的纵横方向）行走的总距离。

2. 特点与适用场景：

• 欧氏距离：

- **特点**：适用于需要测量“直线”距离的场景，反映了两点的真实物理距离。
- **适用场景**：在需要准确反映物理距离的应用中，如物体定位、图像处理中的像素距离计算等。

• 曼哈顿距离：

- **特点**：适用于只能沿着轴方向移动的场景，反映了实际行走的路径距离。
- **适用场景**：在城市街道网格结构中路径规划、在高维空间中用于某些机器学习算法（如K近邻算法）以减少维度影响等。

3. 计算复杂度与性质：

• 欧氏距离：

- 计算涉及平方和开方，可能在高维数据中计算成本较高。
- 满足度量空间的所有性质，包括三角不等式。
- **曼哈顿距离：**
 - 计算简单，仅涉及加法和绝对值运算，适合高维数据的快速计算。
 - 同样满足度量空间的所有性质，但在某些情况下可能不如欧氏距离敏感。

4. 示例比较：

假设有两点A(1,2)和B(4,6)：

- **欧氏距离：**

$$d = \sqrt{(4 - 1)^2 + (6 - 2)^2} = \sqrt{9 + 16} = \sqrt{25} = 5$$

- **曼哈顿距离：**

$$d = |4 - 1| + |6 - 2| = 3 + 4 = 7$$

欧氏距离和曼哈顿距离各有优势，选择哪种距离度量取决于具体应用场景和数据特性。在需要反映真实直线距离的场景中，欧氏距离更为合适；而在需要考虑路径限制或简化计算的场景中，曼哈顿距离则更为适用。

编程实践题

问题：计算雄尾花数据集中各特征之间的Spearman相关系数

1. 介绍

Spearman相关系数是一种非参数统计指标，用于评估两个变量之间的单调关系。与皮尔逊相关系数不同，Spearman相关系数不要求数据服从正态分布，适用于处理非线性关系或序数数据。在雄尾花（Iris）数据集中，计算各特征之间的Spearman相关系数可以帮助我们理解特征之间的依赖关系。

2. 步骤概述

1. **加载数据集：**获取雄尾花数据集，通常可以通过 `scikit-learn` 库直接加载。
2. **数据预处理：**确保数据没有缺失值，并选择需要计算的特征。
3. **计算Spearman相关系数：**使用适当的库函数计算各特征之间的Spearman相关系数。
4. **结果解释：**分析相关系数的大小和方向，理解特征之间的关系。

3. 详细步骤及代码实现

我们将使用Python和相关数据科学库来实现这一过程。

3.1. 导入必要的库

```
import pandas as pd
from scipy.stats import spearmanr
from sklearn.datasets import load_iris
import seaborn as sns
import matplotlib.pyplot as plt
```

- **pandas：**用于数据处理和分析。
- **scipy.stats.spearmanr：**用于计算Spearman相关系数。

- `sklearn.datasets.load_iris`: 加载雄尾花数据集。
- `seaborn` 和 `matplotlib`: 用于可视化相关性矩阵。

3.2. 加载雄尾花数据集

```
# 加载数据
iris = load_iris()
# 创建DataFrame
df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
```

雄尾花数据集包含150个样本，分为三个品种，每个样本有四个特征：萼片长度、萼片宽度、花瓣长度、花瓣宽度。

3.3. 计算Spearman相关系数

```
# 计算Spearman相关系数矩阵
spearman_corr, _ = spearmanr(df)
# 将相关系数转换为DataFrame格式，便于可视化
spearman_corr_df = pd.DataFrame(spearman_corr,
                                index=df.columns,
                                columns=df.columns)
```

- `spearmanr` 函数返回相关系数矩阵和p值矩阵。这里我们只关注相关系数。
- 将相关系数矩阵转换为 `DataFrame`，以便更直观地展示。

3.4. 可视化相关性矩阵

```
# 设置绘图风格
sns.set(style='white')

# 绘制热力图
plt.figure(figsize=(8, 6))
sns.heatmap(spearman_corr_df, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Spearman Correlation Matrix for Iris Dataset')
plt.show()
```

通过热力图，我们可以直观地看到各特征之间的相关性强弱及方向。

3.5. 输出相关系数

```
print("Spearman Correlation Coefficient Matrix:")
print(spearman_corr_df)
```

4. 完整代码示例

```
import pandas as pd
from scipy.stats import spearmanr
from sklearn.datasets import load_iris
import seaborn as sns
import matplotlib.pyplot as plt

# 加载数据
```

```
iris = load_iris()
df = pd.DataFrame(data=iris.data, columns=iris.feature_names)

# 计算Spearman相关系数矩阵
spearman_corr, _ = spearmanr(df)
spearman_corr_df = pd.DataFrame(spearman_corr,
                                index=df.columns,
                                columns=df.columns)

# 输出相关系数矩阵
print("Spearman Correlation Coefficient Matrix:")
print(spearman_corr_df)

# 可视化相关性矩阵
sns.set(style='white')
plt.figure(figsize=(8, 6))
sns.heatmap(spearman_corr_df, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Spearman Correlation Matrix for Iris Dataset')
plt.show()
```

5. 结果解读

运行上述代码后，您将得到一个Spearman相关系数矩阵及其热力图。相关系数范围在-1到1之间：

- 1 表示完全正相关。
- -1 表示完全负相关。
- 0 表示没有相关关系。

通过热力图和矩阵，我们可以观察到哪些特征之间存在强相关性。例如，通常花瓣长度和花瓣宽度之间会有较高的正相关性，这在数据集中应该也能得到验证。

附录1：3-1 correlation.py：

```
import pandas as pd
import matplotlib
matplotlib.use('TkAgg')
from matplotlib import pyplot as plt
data = pd.read_csv('../data/livestreaming.csv', sep=',')
df1 = data.iloc[-24:, :]
# 1. 创建一个画布和两个子图
plt.rcParams.update({'font.size': 20, 'font.sans-serif': ['STSong'], 'axes.unicode_minus': False})
fig, axs = plt.subplots(1, 2, figsize=(14, 4))
# 2. 散点图
axs[0].scatter(data['streamers_count'], data['viewers_count'], color='#F4B183')
axs[0].set(xlabel='主播数量', ylabel='观众数量', title='观众数量和主播数量散点图')
axs[0].ticklabel_format(style='sci', axis='x', scilimits=(0, 0))
# 3. 折线图
df1['time'] = pd.to_datetime(df1['time'])
df1['hour'] = df1['time'].dt.hour
axs[1].plot(df1['hour'], df1['streamers_count'], label='主播数量') # 第一条曲线
axs[1].set(xlabel='时间', ylabel='数量', title='一天内流量变化图')
```

```

axs[1].tick_params(axis='both', which='both', labels=20)
axs[1].ticklabel_format(style='sci', axis='y', scilimits=(0, 0))
ax2 = axs[1].twinx() # 第二条曲线
ax2.plot(df1['hour'], df1['viewers_count'], color='#EC5D3B', label='观众数量')
ax2.set_ylabel('观众数量')
ax2.tick_params(axis='both', which='both', labels=20)
ax2.ticklabel_format(style='sci', axis='y', scilimits=(0, 0))
axs[1].legend(loc='upper left', fontsize=16)
ax2.legend(loc='upper left', bbox_to_anchor=(0, 0.89), fontsize=16)
plt.xticks([0, 3, 7, 11, 15, 19, 23]) # 设置x轴标签格式
plt.subplots_adjust(wspace=0.2) # 调整子图之间的间距
plt.show()
# 4.计算皮尔逊相关系数
pearson_corr = df1['streamers_count'].corr(df1['viewers_count'])

```

附录2：3-2 spearman_kendall.py:

```

import pandas as pd
from scipy.stats import kendalltau, spearmanr
df = pd.read_csv('../data/amazon_reviews.csv') # 读取数据集
df['ordinal_rating'] = df['star_rating'].apply(lambda x: 1 if x <= 2 else 2 if x == 3 else 3) # 将评分转换为数据
# 1.计算kendall相关系数和p值
kendall_corr, kendall_pval = kendalltau(df['ordinal_rating'], df['helpful_votes'])
# 2.计算Spearman相关系数和p值
spearman_corr, spearman_pval = spearmanr(df['ordinal_rating'], df['helpful_votes'])
print('kendall相关系数:', kendall_corr)
print('Spearman相关系数:', spearman_corr)

```

附录3：3-3 nonlinear.py:

```

import pandas as pd
import numpy as np
import matplotlib
matplotlib.use('TkAgg')
import matplotlib.pyplot as plt
from scipy.stats import pearsonr, spearmanr, kendalltau
# from minepy import MINE
from sklearn.feature_selection import mutual_info_regression
# 1.加载数据集
vibration_data = pd.read_csv('../data/bridge.csv')
# 2.确保时间戳被正确解析为datetime对象
vibration_data['timestamp'] = pd.to_datetime(vibration_data['timestamp'])
# 3.可视化振动幅度与时间的关系
plt.figure(figsize=(14, 6))
plt.rcParams['font.sans-serif'] = 'SimSun' # 设置字体为宋体
plt.scatter(vibration_data['timestamp'], vibration_data['amplitude'], alpha=0.6, s=10)
plt.title('桥梁振动幅度与时间关系图', fontsize=20)
plt.xlabel('时间', fontsize=20)
plt.ylabel('震动幅度', fontsize=20)

```

```

plt.grid(False) # 去掉网格
plt.tick_params(axis='both', which='major', labelsize=16) # 设置刻度标签字体大小
plt.show()
# 4.计算相关系数
t_start = vibration_data['timestamp'].min() # 获取数据集中的最早时间
time_in_seconds = (vibration_data['timestamp'] - t_start).dt.total_seconds().values
mutual_info = mutual_info_regression(np.expand_dims(time_in_seconds, axis=1),
vibration_data['amplitude'].values)[0] # 计算互信息
# mine = MINE(alpha=0.6, c=15)
# mine.compute_score(time_in_seconds, vibration_data['amplitude'].values)
# mic_value = mine.mic() # 计算MIC
pearson_corr, pearson_p_value = pearsonr(time_in_seconds,
vibration_data['amplitude'].values) # Pearson相关
spearman_corr, spearman_p_value = spearmanr(time_in_seconds,
vibration_data['amplitude'].values) # Spearman相关
kendall_corr, kendall_p_value = kendalltau(time_in_seconds,
vibration_data['amplitude'].values) # Kendall相关
print("互信息:", mutual_info)
# print("MIC值:", mic_value)
print("Pearson相关系数:", pearson_corr)
print("Spearman相关系数:", spearman_corr)
print("Kendall相关系数:", kendall_corr)

```

附录4：3-4 Partial_Corr.py:

```

import pandas as pd
import seaborn as sns
import matplotlib
matplotlib.use('TkAgg')
import matplotlib.pyplot as plt
from pingouin import partial_corr
# 1. 读取数据集并处理缺失值
data = pd.read_csv('../data/Annual Statistical Indicators Data.csv').dropna()
data.columns = ['地区', '年份'] + ['A'+str(i) for i in range(1, 12)] # 重命名列名
sns.set(font='STSong', font_scale=2.5) # 设置字体和字体缩放
# 2. 计算偏相关系数并绘制热力图
plt.figure(figsize=(24, 12))
# 3. 创建一个空的 DataFrame 来存储偏相关系数
partial_corr_matrix = pd.DataFrame(index=data.columns[2:], columns=data.columns[2:])
# 4. 计算每对变量的偏相关系数
for var1 in data.columns[2:]:
    for var2 in data.columns[2:]:
        if var1 != var2:
            try:
                result = partial_corr(data=data, x=var1, y=var2, covar=[col for col in
data.columns[2:] if col != var1 and col != var2])
                partial_corr_matrix.loc[var1, var2] = result.at['pearson', 'r']
            except Exception as e:
                partial_corr_matrix.loc[var1, var2] = None
        else:
            partial_corr_matrix.loc[var1, var2] = 1.0
partial_corr_matrix = partial_corr_matrix.astype(float)

```



```
plt.subplot(1, 2, 1) # 子图1: 偏相关系数矩阵热力图
sns.heatmap(partial_corr_matrix, cmap='coolwarm', annot=True, fmt='.2f')
plt.title('偏相关系数矩阵热力图')
plt.subplot(1, 2, 2) # 子图2: Pearson相关系数矩阵热力图
corr_matrix = data.iloc[:, 2:].corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.title('Pearson相关系数矩阵热力图')
plt.tight_layout()
plt.show()
```

附录5：3-5 distance.py:

```
import pandas as pd
import matplotlib
matplotlib.use('TkAgg')
from scipy.spatial.distance import euclidean
from scipy.spatial.distance import cosine
# 1.导入数据
data = pd.read_csv('../data/ratings.csv')
# 2.根据userId、movieId、rating三列得到评分行
pivot_df = data.pivot(index='userId', columns='movieId', values='rating')
# 3.填充NaN值为0
pivot_df_subset = pivot_df.fillna(0)
print(pivot_df_subset)
# 4.创建新的DataFrame用于存储矩阵
matrix_df_a = pd.DataFrame(index=pivot_df_subset.index, columns=pivot_df_subset.index)
matrix_df_b = pd.DataFrame(index=pivot_df_subset.index, columns=pivot_df_subset.index)
# 5.填充欧式距离矩阵
for i in pivot_df_subset.index:
    for j in pivot_df_subset.index:
        # 计算用户i和用户j之间的欧式距离
        distance = euclidean(pivot_df_subset.loc[i], pivot_df_subset.loc[j])
        matrix_df_a.loc[i, j] = distance
print(matrix_df_a)
# 6.填充余弦相似度矩阵
for i in pivot_df_subset.index:
    for j in pivot_df_subset.index:
        # 计算用户i和用户j之间的余弦相似度
        similarity = 1 - cosine(pivot_df_subset.loc[i], pivot_df_subset.loc[j])
        matrix_df_b.loc[i, j] = similarity
print(matrix_df_b)
# 7.保存结果
matrix_df_a.to_csv('距离相关分析-欧氏距离.csv')
matrix_df_b.to_csv('距离相关分析-余弦相似度.csv')
```

livestreaming.csv部分数据:

time	streamers_count	viewers_count
2019/3/6 0:00	18981	270126714
2019/3/6 1:00	13892	188527323
2019/3/6 2:00	10587	139727415
2019/3/6 3:00	7890	104173437
2019/3/6 4:00	6487	82368147
2019/3/6 5:00	5443	65856669
2019/3/6 6:00	4853	59062388
2019/3/6 7:00	4773	59871640
2019/3/6 8:00	5858	71374021
2019/3/6 9:00	7830	93140327
2019/3/6 10:00	9911	115152445
2019/3/6 11:00	11267	128426750
2019/3/6 12:00	13071	148432417
2019/3/6 13:00	17612	162476040
2019/3/6 14:00	22289	177804128
2019/3/6 15:00	23680	193932883
2019/3/6 16:00	24073	200714670
2019/3/6 17:00	23983	205169204
2019/3/6 18:00	22181	221588024
2019/3/6 19:00	25617	276252703
2019/3/6 20:00	30985	327783240
2019/3/6 21:00	33193	363431743
2019/3/6 22:00	32928	352375781
2019/3/6 23:00	27238	331629355
2019/3/7 0:00	19644	275985468
2019/3/7 1:00	14474	199045299
2019/3/7 2:00	11113	150084503
2019/3/7 3:00	8243	108723097

ratings.csv:

userId	movieId	rating	timestamp
1	1	4	964982703
1	3	4	964981247
1	6	4	964982224
1	47	5	964983815
1	50	5	964982931
1	70	3	964982400
1	101	5	964980868
1	110	4	964982176
1	151	5	964984041
1	157	5	964984100
1	163	5	964983650
1	216	5	964981208
1	223	3	964980985
1	231	5	964981179
1	235	4	964980908
1	260	5	964981680
1	296	3	964982967
1	316	3	964982310
1	333	5	964981179
1	349	4	964982563
1	356	4	964980962
1	362	5	964982588
1	367	4	964981710
1	423	3	964982363
1	441	4	964980868
1	457	5	964981909
1	480	4	964982346
1	500	3	964981208

Annual Statistical Indictors Data.csv:

地区	年份	国内生产总值	第一产业增加值	第二产业增加值	第三产业增加值	社会商品零售总额	货物进出口总额	年末总人口	在岗职工平均工资	普通高等学校在校学生数	医院、卫生院数	房地产开发投资额
北京	2016	25669.13	129.79	4944.44	20594.9	11005.1	282348.96	1362.86	122749	59.9188	713	4000.57
天津	2016	17885.39	220.22	7571.35	10093.82	5635.8	102655.95	1044.4	87806	51.3842	571	2300.01
石家庄	2016	5927.73	480.88	2693.92	2752.93	2975.2	12160.29	1038	61189	44.1812	425	1015.77
太原	2016	2955.6	38.77	1067.49	1849.34	1666.2	13314.33	370.25	64820	43.2234	245	680.13
呼和浩特	2016	3173.59	113.49	884.43	2175.67	1481.5	1308.67	240.97	56213	23.7734	185	520.52
沈阳	2016	5460.01	266.36	2135.63	3058.02	3985.9	11331.41	734.4	67444	40.3589	383	709.67
大连	2016	6730.33	462.78	2793.69	3473.86	3410.1	51443.88	595.63	73764	29.0217	315	535.17
长春	2016	5917.94	323.53	2915.66	2678.74	2650.3	14162.91	753.43	68434	43.4366	297	596.65
哈尔滨	2016	6101.61	691.18	1896.66	3513.77	3744.2	3974.63	962.05	62583	63.624	471	526.13
上海	2016	28178.65	109.47	8406.28	19662.9	10946.6	433768.19	1450	120503	51.4683	656	3709.03
南京	2016	10503.02	252.54	4117.32	6133.16	5088.2	50214.06	662.79	90191	82.7773	225	1845.6
杭州	2016	11313.72	304.21	4120.93	6888.59	5176.2	67992.41	736	87153	42.7978	365	2606.63
宁波	2016	8686.49	302.06	4455.34	3929.1	3667.6	94923.22	590.96	83656	15.5144	251	1270.33
合肥	2016	6274.38	270.17	3181.24	2822.97	2445.7	18686.99	729.83	71054	49.9515	462	1352.59
福州	2016	6197.64	492.25	2590.43	3114.96	3763.1	31978.54	687.06	67630	31.7477	230	1679.44
厦门	2016	3784.27	23.19	1544.59	2216.49	1283.5	77176.81	220.55	69218	14.2948	60	765.8
南昌	2016	4354.99	181.77	2307.24	1865.98	1868	9407.28	522.79	65812	61.1819	194	674.6
济南	2016	6536.12	317.31	2368.9	3849.91	3764.8	9830.7	632.83	77012	72.6301	270	1164.14
青岛	2016	10011.29	371.01	4160.67	5479.61	4104.9	65581.15	791.35	76616	34.0875	322	1369.14
郑州	2016	8025.31	156.35	3728.66	4140.29	3665.8	55028.78	827.06	61149	88.9329	317	2778.95
武汉	2016	11912.61	390.62	5227.05	6294.94	5610.6	23780.93	833.85	71963	94.8768	386	2517.44
长沙	2016	9455.36	370.95	4521.02	4563.4	4117.4	10934.58	696	77782	59.002	286	1266.63
广州	2016	19547.44	239.28	5751.59	13556.57	8706.5	129308.95	870.49	89096	105.7281	273	2540.85

amazon_reviews.csv:

marketplace	product_id	review_id	helpful_votes	star_rating
US	B00L9EPT8O	R3OW6KZQVW8JG	0	1
US	B00L9EPT8O	RVN4P3GJ8QMTA	1	2
US	B00L9EPT8O	R4W1A8Z4LJH7D	2	3
US	B00L9EPT8O	R3K3W3YQYJDP6U	3	4
US	B00L9EPT8O	R1JZV7R5EJX8OW	4	5
US	B00L9EPT8O	R2U5MIYOYUT9KU	0	1
US	B00L9EPT8O	R1KZ9XGJ9I9JSN	1	2
US	B00L9EPT8O	R3W0C3O3J4Z4J4	2	3
US	B00L9EPT8O	R3P4Z4P4Z4P4Z4	3	4
US	B00L9EPT8O	R2Z4P4Z4P4Z4P4	4	5
US	B00L9EPT8O	R5A5P5Z5P5Z5P5	0	1
US	B00L9EPT8O	R3Z3P3Z3P3Z3P3	1	2
US	B00L9EPT8O	R2W2P2Z2P2Z2P2	2	3
US	B00L9EPT8O	R1U1P1Z1P1Z1P1	3	4
US	B00L9EPT8O	R3K3P3Z3P3Z3P3	4	5
US	B00L9EPT8O	R4W4P4Z4P4Z4P4	0	1
US	B00L9EPT8O	R2K2P2Z2P2Z2P2	1	2
US	B00L9EPT8O	R1J1P1Z1P1Z1P1	2	3
US	B00L9EPT8O	R3W3P3Z3P3Z3P3	3	4
US	B00L9EPT8O	R4U4P4Z4P4Z4P4	4	5
US	B00L9EPT8O	R5W5P5Z5P5Z5P5	0	1
US	B00L9EPT8O	R3Z3P3Z3P3Z3P3	1	2
US	B00L9EPT8O	R2W2P2Z2P2Z2P2	2	3
US	B00L9EPT8O	R1U1P1Z1P1Z1P1	3	4
US	B00L9EPT8O	R3K3P3Z3P3Z3P3	4	5

bridge.csv:

timestamp	amplitude	frequency
2020/8/6 0:00	15.29215704	1.031876124
2020/8/6 0:01	11.27592735	0.90983788
2020/8/6 0:03	13.08710822	0.973720989
2020/8/6 0:05	16.94897804	1.064459506
2020/8/6 0:06	15.90432504	0.850834456
2020/8/6 0:08	7.445101364	0.972911832
2020/8/6 0:10	13.30239834	0.966649686
2020/8/6 0:12	10.07315657	1.017040408
2020/8/6 0:13	10.29254666	1.008765472
2020/8/6 0:15	11.90883659	0.947090674
2020/8/6 0:17	11.18385547	1.105093452
2020/8/6 0:19	15.18905774	0.826891616
2020/8/6 0:20	13.18367468	0.814215389
2020/8/6 0:22	11.33970572	0.968103736
2020/8/6 0:24	12.38016756	1.1586777
2020/8/6 0:25	12.12325921	0.91102774
2020/8/6 0:27	15.67787621	0.86892373
2020/8/6 0:29	10.65329466	1.171283358
2020/8/6 0:31	12.28081404	0.867289302
2020/8/6 0:32	8.851859643	0.949936215
2020/8/6 0:34	3.827391261	1.091606204
2020/8/6 0:36	13.51909181	0.904369603
2020/8/6 0:38	14.22306505	0.891939811
2020/8/6 0:39	9.474410626	0.959066894
2020/8/6 0:41	18.58093141	1.114945926
2020/8/6 0:43	7.478928848	1.164393341